

# Revisiting the Nyström Method for Improved Large-scale Machine Learning

**Alex Gittens**

GITTENS@ICSI.BERKELEY.EDU

**Michael W. Mahoney**

MMAHONEY@STAT.BERKELEY.EDU

*International Computer Science Institute and Department of Statistics*

*University of California, Berkeley*

*Berkeley, CA*

**Editor:** Mehryar Mohri

## Abstract

We reconsider randomized algorithms for the low-rank approximation of symmetric positive semi-definite (SPSD) matrices such as Laplacian and kernel matrices that arise in data analysis and machine learning applications. Our main results consist of an empirical evaluation of the performance quality and running time of sampling and projection methods on a diverse suite of SPSPD matrices. Our results highlight complementary aspects of sampling versus projection methods; they characterize the effects of common data preprocessing steps on the performance of these algorithms; and they point to important differences between uniform sampling and nonuniform sampling methods based on leverage scores. In addition, our empirical results illustrate that existing theory is so weak that it does not provide even a qualitative guide to practice. Thus, we complement our empirical results with a suite of worst-case theoretical bounds for both random sampling and random projection methods. These bounds are qualitatively superior to existing bounds—*e.g.*, improved additive-error bounds for spectral and Frobenius norm error and relative-error bounds for trace norm error—and they point to future directions to make these algorithms useful in even larger-scale machine learning applications.

**Keywords:** Nyström approximation, low-rank approximation, kernel methods, randomized algorithms, numerical linear algebra

## 1. Introduction

We reconsider randomized algorithms for the low-rank approximation of symmetric positive semi-definite (SPSD) matrices such as Laplacian and kernel matrices that arise in data analysis and machine learning applications. Our goal is to obtain an improved understanding, both empirically and theoretically, of the complementary strengths of sampling versus projection methods on realistic data. Our main results consist of an empirical evaluation of the performance quality and running time of sampling and projection methods on a diverse suite of dense and sparse SPSPD matrices drawn both from machine learning as well as more general data analysis applications. These results are not intended to be comprehensive but instead to be illustrative of how randomized algorithms for the low-rank approximation of SPSPD matrices behave in a broad range of realistic machine learning and data analysis applications.

Our empirical results point to several directions that are not explained well by existing theory. (For example, that the results are much better than existing worst-case theory would suggest, and that sampling with respect to the statistical leverage scores leads to results that are complementary to those achieved by projection-based methods.) Thus, we complement our empirical results with a suite of worst-case theoretical bounds for both random sampling and random projection methods. These bounds are qualitatively superior to existing bounds—*e.g.*, improved additive-error bounds for spectral and Frobenius norm error and relative-error bounds for trace norm error. By considering random sampling and random projection algorithms on an equal footing, we identify within our analysis deterministic structural properties of the input data and sampling/projection methods that are responsible for high-quality low-rank approximation.

In more detail, our main contributions are fourfold.

- First, we provide an empirical illustration of the complementary strengths and weaknesses of data-independent random projection methods and data-dependent random sampling methods when applied to SPSD matrices. We do so for a diverse class of SPSD matrices drawn from machine learning and data analysis applications, and we consider reconstruction error with respect to the spectral, Frobenius, and trace norms. Depending on the parameter settings, the matrix norm of interest, the data set under consideration, etc., one or the other method might be preferable. In addition, we illustrate how these empirical properties can often be understood in terms of the structural nonuniformities of the input data that are of independent interest.
- Second, we consider the running time of high-quality sampling and projection algorithms. For random sampling algorithms, the computational bottleneck is typically the exact or approximate computation of the importance sampling distribution with respect to which one samples; and for random projection methods, the computational bottleneck is often the implementation of the random projection. By exploiting and extending recent work on “fast” random projections and related recent work on “fast” approximation of the statistical leverage scores, we illustrate that high-quality leverage-based random sampling and high-quality random projection algorithms have comparable running times. Although both are slower than simple (and in general much lower-quality) uniform sampling, both can be implemented more quickly than a naïve computation of an orthogonal basis for the top part of the spectrum.
- Third, our main technical contribution is a set of deterministic structural results that hold for any “sketching matrix” applied to an SPSD matrix. We call these “deterministic structural results” since there is no randomness involved in their statement or analysis and since they depend on structural properties of the input data matrix and the way the sketching matrix interacts with the input data. In particular, they highlight the importance of the statistical leverage scores, which have proven important in other applications of random sampling and random projection algorithms.
- Fourth, our main algorithmic contribution is to show that when the low-rank sketching matrix represents certain random projection or random sampling operations, then we obtain worst-case quality-of-approximation bounds that hold with high probability. These bounds are qualitatively better than existing bounds and they illustrate

how high-quality random sampling algorithms and high-quality random projection algorithms can be treated from a unified perspective.

A novel aspect of our work is that we adopt a unified approach to these low-rank approximation questions—unified in the sense that we consider both sampling and projection algorithms on an equal footing, and that we illustrate how the structural nonuniformities responsible for high-quality low-rank approximation in worst-case analysis also have important empirical consequences in a diverse class of SPSD matrices. By identifying deterministic structural conditions responsible for high-quality low-rank approximation of SPSD matrices, we highlight complementary aspects of sampling and projection methods; and by illustrating the empirical consequences of structural nonuniformities, we provide theory that is a much closer guide to practice than has been provided by prior work. We note also that our deterministic structural results could be used to check, in an *a posteriori* manner, the quality of a sketching method for which one cannot establish an *a priori* bound.

Our analysis is timely for several reasons. First, in spite of the empirical successes of Nyström-based and other randomized low-rank methods, existing theory for the Nyström method is quite modest. For example, existing worst-case bounds such as those of Drineas and Mahoney (2005) are very weak, especially compared with existing bounds for least-squares regression and general low-rank matrix approximation problems (Drineas et al., 2008, 2010; Mahoney, 2011).<sup>1</sup> Moreover, many other worst-case bounds make very strong assumptions about the coherence properties of the input data (Kumar et al., 2012; Gittens, 2012). Second, there have been conflicting views in the literature about the usefulness of uniform sampling versus nonuniform sampling based on the empirical statistical leverage scores of the data in realistic data analysis and machine learning applications. For example, some work has concluded that the statistical leverage scores of realistic data matrices are fairly uniform, meaning that the coherence is small and thus uniform sampling is appropriate (Williams and Seeger, 2001; Kumar et al., 2012); while other work has demonstrated that leverage scores are often very nonuniform in ways that render uniform sampling inappropriate and that can be essential to highlight properties of downstream interest (Paschou et al., 2007; Mahoney and Drineas, 2009). Third, in recent years several high-quality numerical implementations of randomized matrix algorithms for least-squares and low-rank approximation problems have been developed (Avron et al., 2010; Meng et al., 2014; Woolfe et al., 2008; Rokhlin et al., 2009; Martinsson et al., 2011). These have been developed from a “scientific computing” perspective, where condition numbers, spectral norms, etc. are of greater interest (Mahoney, 2012), and where relatively strong homogeneity assumptions can be made about the input data. In many “data analytics” applications, the questions one asks are very different, and the input data are much less well-structured. Thus, we expect that some of our results will help guide the development of algorithms and implementations that are more appropriate for large-scale analytics applications.

In the next section, Section 2, we start by presenting some notation, preliminaries, and related prior work. Then, in Section 3 we present our main empirical results; and in

---

1. This statement may at first surprise the reader, since an SPSD matrix is an example of a general matrix, and one might suppose that the existing theory for general matrices could be applied to SPSD matrices. While this is true, these existing methods for general matrices do not in general respect the symmetry or positive semi-definiteness of the input.

Section 4 we present our main theoretical results. We conclude in Section 5 with a brief discussion of our results in a broader context.

## 2. Notation, Preliminaries, and Related Prior Work

In this section, we introduce the notation used throughout the paper, and we address several preliminary considerations, including reviewing related prior work.

### 2.1 Notation

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be an arbitrary SPSD matrix with eigenvalue decomposition  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T$ , where we partition  $\mathbf{U}$  and  $\mathbf{\Sigma}$  as

$$\mathbf{U} = (\mathbf{U}_1 \quad \mathbf{U}_2) \quad \text{and} \quad \mathbf{\Sigma} = \begin{pmatrix} \mathbf{\Sigma}_1 & \\ & \mathbf{\Sigma}_2 \end{pmatrix}. \quad (1)$$

Here,  $\mathbf{U}_1$  has  $k$  columns and spans the top  $k$ -dimensional eigenspace of  $\mathbf{A}$ , and  $\mathbf{\Sigma}_1 \in \mathbb{R}^{k \times k}$  is full-rank.<sup>2</sup> We denote the eigenvalues of  $\mathbf{A}$  with  $\lambda_1(\mathbf{A}) \geq \dots \geq \lambda_n(\mathbf{A})$ .

Given  $\mathbf{A}$  and a rank parameter  $k$ , the *statistical leverage scores of  $\mathbf{A}$  relative to the best rank- $k$  approximation to  $\mathbf{A}$*  equal the squared Euclidean norms of the rows of the  $n \times k$  matrix  $\mathbf{U}_1$ :

$$\ell_j = \|(\mathbf{U}_1)_j\|^2. \quad (2)$$

The leverage scores provide a more refined notion of the structural nonuniformities of  $\mathbf{A}$  than does the notion of *coherence*,  $\mu = \frac{n}{k} \max_{i \in \{1, \dots, n\}} \ell_i$ , which equals (up to scale) the largest leverage score; and they have been used historically in regression diagnostics to identify particularly influential or outlying data points. Less obviously, the statistical leverage scores play a crucial role in recent work on randomized matrix algorithms: they define the key structural nonuniformity that must be dealt with in order to obtain high-quality low-rank and least-squares approximation of general matrices via random sampling and random projection methods (Mahoney, 2011). Although Equation (2) defines them with respect to a particular basis, the statistical leverage scores equal the diagonal elements of the projection matrix onto the span of that basis, and thus they can be computed from any basis spanning the same space. Moreover, they can be approximated more quickly than the time required to compute that basis with a truncated SVD or a QR decomposition (Drineas et al., 2012).

We denote by  $\mathbf{S}$  an arbitrary  $n \times \ell$  “sketching” matrix that, when post-multiplying a matrix  $\mathbf{A}$ , maps points from  $\mathbb{R}^n$  to  $\mathbb{R}^\ell$ . We are most interested in the case where  $\mathbf{S}$  is a random matrix that represents a random sampling process or a random projection process, but we do not impose this as a restriction unless explicitly stated. We let

$$\mathbf{\Omega}_1 = \mathbf{U}_1^T \mathbf{S} \quad \text{and} \quad \mathbf{\Omega}_2 = \mathbf{U}_2^T \mathbf{S} \quad (3)$$

denote the projection of  $\mathbf{S}$  onto the top and bottom eigenspaces of  $\mathbf{A}$ , respectively.

---

2. Variants of our results hold trivially if the rank of  $\mathbf{A}$  is  $k$  or less, so we focus on this more general case here.

Recall that, by keeping just the top  $k$  singular vectors, the matrix  $\mathbf{A}_k := \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{U}_1^T$  is the best rank- $k$  approximation to  $\mathbf{A}$ , when measured with respect to any unitarily-invariant matrix norm, *e.g.*, the spectral, Frobenius, or trace norm. For a vector  $\mathbf{x} \in \mathbb{R}^n$ , let  $\|\mathbf{x}\|_\xi$ , for  $\xi = 1, 2, \infty$ , denote the 1-norm, the Euclidean norm, and the  $\infty$ -norm, respectively, and let  $\text{Diag}(\mathbf{A})$  denote the vector consisting of the diagonal entries of the matrix  $\mathbf{A}$ . Then,  $\|\mathbf{A}\|_2 = \|\text{Diag}(\boldsymbol{\Sigma})\|_\infty$  denotes the *spectral norm* of  $\mathbf{A}$ ;  $\|\mathbf{A}\|_F = \|\text{Diag}(\boldsymbol{\Sigma})\|_2$  denotes the *Frobenius norm* of  $\mathbf{A}$ ; and  $\|\mathbf{A}\|_\star = \|\text{Diag}(\boldsymbol{\Sigma})\|_1$  denotes the *trace norm* (or nuclear norm) of  $\mathbf{A}$ . Clearly,

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \|\mathbf{A}\|_\star \leq \sqrt{n} \|\mathbf{A}\|_F \leq n \|\mathbf{A}\|_2.$$

We quantify the quality of our algorithms by the “additional error” (above and beyond that incurred by the best rank- $k$  approximation to  $\mathbf{A}$ ). In the theory of algorithms, bounds of the form provided by (16) below are known as *additive-error bounds*, the reason being that the additional error is an additive factor of the form  $\epsilon$  times a size scale that is larger than the “base error” incurred by the best rank- $k$  approximation. In this case, the goal is to minimize the “size scale” of the additional error. Bounds of this form are very different and in general weaker than when the additional error enters as a multiplicative factor, such as when the error bounds are of the form  $\|\mathbf{A} - \tilde{\mathbf{A}}\| \leq f(n, k, \eta) \|\mathbf{A} - \mathbf{A}_k\|$ , where  $f(\cdot)$  is some function and  $\eta$  represents other parameters of the problem. These latter bounds are of greatest interest when  $f = 1 + \epsilon$ , for an error parameter  $\epsilon$ , as in (18) and (19) below. These *relative-error bounds*, in which the size scale of the additional error equals that of the base error, provide a *much* stronger notion of approximation than additive-error bounds.

## 2.2 Preliminaries

In many machine learning and data analysis applications, one is interested in symmetric positive semi-definite (SPSD) matrices, *e.g.*, kernel matrices and Laplacian matrices. One common column-sampling-based approach to low-rank approximation of SPSPD matrices is the so-called Nyström method (Williams and Seeger, 2001; Drineas and Mahoney, 2005; Kumar et al., 2012). The Nyström method—both randomized and deterministic variants—has proven useful in applications where the kernel matrices are reasonably well-approximated by low-rank matrices; and it has been applied to Gaussian process regression, spectral clustering and image segmentation, manifold learning, and a range of other common machine learning tasks (Williams and Seeger, 2001; Williams et al., 2002; Fowlkes et al., 2004; Talwalkar et al., 2008; Zhang and Kwok, 2010; Kumar et al., 2012). The simplest Nyström-based procedure selects columns from the original data set uniformly at random and then uses those columns to construct a low-rank SPSPD approximation. Although this procedure can be effective in practice for certain input matrices, two extensions (both of which are more expensive) can substantially improve the performance, *e.g.*, lead to lower reconstruction error for a fixed number of column samples, both in theory and in practice. The first extension is to sample columns with a judiciously-chosen nonuniform importance sampling distribution; and the second extension is to randomly mix (or combine linearly) columns before sampling them. For the random sampling algorithms, an important question is what importance sampling distribution should be used to construct the sample; while for the random projection algorithms, an important question is how to implement the random projections. In either case, appropriate consideration should be paid to questions such as

whether the data are sparse or dense, how the eigenvalue spectrum decays, the nonuniformity properties of eigenvectors, *e.g.*, as quantified by the statistical leverage scores, whether one is interested in reconstructing the matrix or performing a downstream machine learning task, and so on.

The following sketching model subsumes both of these classes of methods.

- *SPSD Sketching Model.* Let  $\mathbf{A}$  be an  $n \times n$  positive semi-definite matrix, and let  $\mathbf{S}$  be a matrix of size  $n \times \ell$ , where  $\ell \ll n$ . Take

$$\mathbf{C} = \mathbf{A}\mathbf{S} \quad \text{and} \quad \mathbf{W} = \mathbf{S}^T \mathbf{A}\mathbf{S}.$$

Then  $\mathbf{C}\mathbf{W}^\dagger \mathbf{C}^T$  is a low-rank approximation to  $\mathbf{A}$  with rank at most  $\ell$ .

We should note that the SPSPD Sketching Model, formulated in this way, is *not* guaranteed to be numerically stable: if  $\mathbf{W}$  is ill-conditioned, then instabilities may arise in forming the product  $\mathbf{C}\mathbf{W}^\dagger \mathbf{C}^T$ . For simplicity in our presentation, we do not describe the generalizations of our results that could be obtained for the various algorithmic tweaks that have been considered to address this potential issue (Drineas et al., 2008; Mahoney and Drineas, 2009; Chiu and Demanet, 2013).

The choice of distribution for the sketching matrix  $\mathbf{S}$  leads to different classes of low-rank approximations. For example, if  $\mathbf{S}$  represents the process of column sampling, either uniformly or according to a nonuniform importance sampling distribution, then we refer to the resulting approximation as a Nyström extension; if  $\mathbf{S}$  consists of random linear combinations of most or all of the columns of  $\mathbf{A}$ , then we refer to the resulting approximation as a projection-based SPSPD approximation. In this paper, we focus on Nyström extensions and projection-based SPSPD approximations that fit the above SPSPD Sketching Model. In particular, we do not consider adaptive schemes, which iteratively select columns to progressively decrease the approximation error. While these methods often perform well in practice (Belabbas and Wolfe, 2009b,a; Farahat et al., 2011; Kumar et al., 2012), rigorous analyses of them are hard to come by—interested readers are referred to the discussion in (Farahat et al., 2011; Kumar et al., 2012).

### 2.3 The Power Method

One can obtain the optimal rank- $k$  approximation to  $\mathbf{A}$  by forming an SPSPD sketch where the sketching matrix  $\mathbf{S}$  is an orthonormal basis for the range of  $\mathbf{A}_k$ , because with such a choice,

$$\mathbf{C}\mathbf{W}^\dagger \mathbf{C}^T = \mathbf{A}\mathbf{S}(\mathbf{S}^T \mathbf{A}\mathbf{S})^\dagger \mathbf{S}^T \mathbf{A} = \mathbf{A}(\mathbf{S}\mathbf{S}^T \mathbf{A}\mathbf{S}\mathbf{S}^T)^\dagger \mathbf{A} = \mathbf{A}(\mathbf{P}_{\mathbf{A}_k} \mathbf{A} \mathbf{P}_{\mathbf{A}_k})^\dagger \mathbf{A} = \mathbf{A}\mathbf{A}_k^\dagger \mathbf{A} = \mathbf{A}_k.$$

Of course, one cannot quickly obtain such a basis; this motivates considering sketching matrices  $\mathbf{S}_q$  obtained using the power method: that is, taking  $\mathbf{S}_q = \mathbf{A}^q \mathbf{S}_0$  where  $q$  is a positive integer and  $\mathbf{S}_0 \in \mathbb{R}^{n \times \ell}$  with  $l \geq k$ . As  $q \rightarrow \infty$ , assuming  $\mathbf{U}_1^T \mathbf{S}_0$  has full row-rank, the matrices  $\mathbf{S}_q$  increasingly capture the dominant  $k$ -dimensional eigenspaces of  $\mathbf{A}$  (see Golub and Van Loan, 1996, Chapter 8), so one can reasonably expect that the sketching matrix  $\mathbf{S}_q$  produces SPSPD sketches of  $\mathbf{A}$  with lower additional error.

SPSPD sketches produced using  $q$  iterations of the power method have lower error than sketches produced without using the power method, but are roughly  $q$  times more costly to

produce. Thus, the power method is most applicable when  $\mathbf{A}$  is such that one can compute the product  $\mathbf{A}^q \mathbf{S}_0$  fast. We consider the empirical performance of sketches produced using the power method in Section 3, and we consider the theoretical performance in Section 4.

## 2.4 Related Prior Work

Motivated by large-scale data analysis and machine learning applications, recent theoretical and empirical work has focused on “sketching” methods such as random sampling and random projection algorithms. A large part of the recent body of this work on randomized matrix algorithms has been summarized in the recent monograph by Mahoney (2011) and the recent review article by Halko et al. (2011). Here, we note that, on the empirical side, both random projection methods (*e.g.*, Bingham and Mannila, 2001; Fradkin and Madigan, 2003; Venkatasubramanian and Wang, 2011; Banerjee et al., 2012) and random sampling methods (*e.g.*, Paschou et al., 2007; Mahoney and Drineas, 2009) have been used in applications for clustering and classification of general data matrices; and that some of this work has highlighted the importance of the statistical leverage scores that we use in this paper (Paschou et al., 2007; Mahoney and Drineas, 2009; Mahoney, 2011; Yip et al., 2014). In parallel, so-called Nyström-based methods have also been used in machine learning applications. Originally used by Williams and Seeger to solve regression and classification problems involving Gaussian processes when the SPSP matrix  $\mathbf{A}$  is well-approximated by a low-rank matrix (Williams and Seeger, 2001; Williams et al., 2002), the Nyström extension has been used in a large body of subsequent work. For example, applications of the Nyström method to large-scale machine learning problems include the work of Talwalkar et al. (2008); Kumar et al. (2009a,c); Mackey et al. (2011b) and Zhang et al. (2008); Li et al. (2010); Zhang and Kwok (2010), and applications in statistics and signal processing include the work of Parker et al. (2005); Belabbas and Wolfe (2007a,b); Spendley and Wolfe (2008); Belabbas and Wolfe (2008, 2009b,a).

Much of this work has focused on new proposals for selecting columns (*e.g.*, Zhang et al., 2008; Zhang and Kwok, 2009; Liu et al., 2010; Arcolano and Wolfe, 2010; Li et al., 2010) and/or coupling the method with downstream applications (*e.g.*, Bach and Jordan, 2005; Cortes et al., 2010; Jin et al., 2013; Homrighausen and McDonald, 2011; Machart et al., 2011; Bach, 2013). The most detailed results are provided by Kumar et al. (2012) as well as the conference papers on which it is based (Kumar et al., 2009a,b,c). Interestingly, they observe that uniform sampling performs quite well, suggesting that in the data they considered the leverage scores are quite uniform, which also motivated the related works of Talwalkar and Rostamizadeh (2010); Mohri and Talwalkar (2011). This is in contrast with applications in genetics (Paschou et al., 2007), term-document analysis (Mahoney and Drineas, 2009), and astronomy (Yip et al., 2014), where the statistical leverage scores were seen to be very nonuniform in ways of interest to the downstream scientist; we return to this issue in Section 3.

On the theoretical side, much of the work has followed that of Drineas and Mahoney (2005), who provided the first rigorous bounds for the Nyström extension of a general SPSP matrix. They show that when  $\Omega(k\epsilon^{-4} \ln \delta^{-1})$  columns are sampled with an importance sampling distribution that is proportional to the square of the diagonal entries of  $\mathbf{A}$ , then

$$\|\mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T\|_\xi \leq \|\mathbf{A} - \mathbf{A}_k\|_\xi + \epsilon \sum_{k=1}^n (\mathbf{A})_{ii}^2 \quad (4)$$

holds with probability  $1 - \delta$ , where  $\xi = 2, F$  represents the Frobenius or spectral norm. (Actually, they prove a stronger result of the form given in Equation (4), except with  $\mathbf{W}^\dagger$  replaced with  $\mathbf{W}_k^\dagger$ , where  $\mathbf{W}_k$  represents the best rank- $k$  approximation to  $\mathbf{W}$  (Drineas and Mahoney, 2005).) Subsequently, Kumar, Mohri, and Talwalkar show that if  $\mu k \ln(k/\delta)$  columns are sampled uniformly at random with replacement from an  $\mathbf{A}$  that has *exactly* rank  $k$ , then one achieves exact recovery, *i.e.*,  $\mathbf{A} = \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T$ , with high probability (Kumar et al., 2009a). Gittens (2012) extends this to the case where  $\mathbf{A}$  is only approximately low-rank. In particular, he shows that if  $\ell = \Omega(\mu k \ln k)$  columns are sampled uniformly at random (either with or without replacement), then

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 \leq \|\mathbf{A} - \mathbf{A}_k\|_2 \left( 1 + \frac{2n}{\ell} \right) \quad (5)$$

with probability exceeding  $1 - \delta$  and

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 \leq \|\mathbf{A} - \mathbf{A}_k\|_2 + \frac{2}{\delta} \cdot \|\mathbf{A} - \mathbf{A}_k\|_* \quad (6)$$

with probability exceeding  $1 - 2\delta$ .

We have described these prior theoretical bounds in detail to emphasize how strong, relative to the prior work, our new bounds are. For example, Equation (4) provides an additive-error approximation with a very large scale; the bounds of Kumar, Mohri, and Talwalkar require a sampling complexity that depends on the coherence of the input matrix (Kumar et al., 2009a), which means that unless the coherence is very low one needs to sample essentially all the rows and columns in order to reconstruct the matrix; Equation (5) provides a bound where the additive scale depends on  $n$ ; and Equation (6) provides a spectral norm bound where the scale of the additional error is the (much larger) trace norm. Table 1 compares the bounds on the approximation errors of SPSD sketches derived in this work to those available in the literature. We note further that Wang and Zhang recently established lower-bounds on the worst-case relative spectral and trace norm errors of uniform Nyström extensions (Wang and Zhang, 2013). Our Lemma 8 provides matching upper bounds, showing the optimality of these estimates.

A related stream of research concerns projection-based low-rank approximations of general (*i.e.*, non-SPSD) matrices (Halko et al., 2011; Mahoney, 2011). Such approximations are formed by first constructing an approximate basis for the top left invariant subspace of  $\mathbf{A}$ , and then restricting  $\mathbf{A}$  to this space. Algorithmically, one constructs  $\mathbf{Y} = \mathbf{A}\mathbf{S}$ , where  $\mathbf{S}$  is a sketching matrix, then takes  $\mathbf{Q}$  to be a basis obtained from the QR decomposition of  $\mathbf{Y}$ , and then forms the low-rank approximation  $\mathbf{Q}\mathbf{Q}^T\mathbf{A}$ . The survey paper Halko et al. (2011) proposes two schemes for the approximation of SPSD matrices that fit within this paradigm:  $\mathbf{Q}(\mathbf{Q}^T\mathbf{A}\mathbf{Q})\mathbf{Q}^T$  and  $(\mathbf{A}\mathbf{Q})(\mathbf{Q}^T\mathbf{A}\mathbf{Q})^\dagger(\mathbf{Q}^T\mathbf{A})$ . The first scheme—for which Halko et al. (2011) provides quite sharp error bounds when  $\mathbf{S}$  is a matrix of i.i.d. standard Gaussian random variables—has the salutary property of being numerically stable. In Wang and Zhang (2013), the authors show that using the first scheme with an adaptively sampled  $\mathbf{S}$  results in approximations with expected Frobenius error within a factor of  $1 + \epsilon$  of the optimal rank- $k$  approximation error when  $O(k/\epsilon^2)$  columns are sampled.

Halko et al. (2011) does not provide any theoretical guarantees for the second scheme, but observes that this latter scheme produces noticeably more accurate approximations in



Source	$\ell$	$\ \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T\ _2$	$\ \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T\ _F$	$\ \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T\ _\star$
Prior works				
Drineas and Mahoney (2005)	$\Omega(\epsilon^{-4}k)$	$\text{opt}_2 + \epsilon \sum_{i=1}^n A_{ii}^2$	$\text{opt}_F + \epsilon \sum_{i=1}^n A_{ii}^2$	–
Belabbas and Wolfe (2009b)	$\Omega(1)$	–	–	$O\left(\frac{n-\ell}{n}\right) \ \mathbf{A}\ _\star$
Talwalkar and Rostamizadeh (2010)	$\Omega(\mu_r r \ln r)$	0	0	0
Kumar et al. (2012)	$\Omega(1)$	$\text{opt}_2 + \frac{n}{\sqrt{\ell}} \ \mathbf{A}\ _2$	$\text{opt}_F + n\left(\frac{k}{\ell}\right)^{1/4} \ \mathbf{A}\ _2$	–
This work				
Lemma 8, uniform column sampling	$\Omega\left(\frac{\mu_k k \ln k}{(1-\epsilon)^2}\right)$	$\text{opt}_2(1 + \frac{n}{\epsilon\ell})$	$\text{opt}_F + \epsilon^{-1}\text{opt}_\star$	$\text{opt}_\star(1 + \epsilon^{-1})$
Lemma 5, leverage-based column sampling	$\Omega\left(\frac{k \ln(k/\beta)}{\beta\epsilon^2}\right)$	$\text{opt}_2 + \epsilon^2\text{opt}_\star$	$\text{opt}_F + \epsilon\text{opt}_\star$	$(1 + \epsilon^2)\text{opt}_\star$
Lemma 6, Fourier-based projection	$\Omega(\epsilon^{-1}k \ln n)$	$(1 + \frac{1}{1-\sqrt{\epsilon}})\text{opt}_2 + \frac{\epsilon\text{opt}_\star}{(1-\sqrt{\epsilon})k}$	$\text{opt}_F + \sqrt{\epsilon}\text{opt}_\star$	$(1 + \epsilon)\text{opt}_\star$
Lemma 7, Gaussian-based projection	$\Omega(k\epsilon^{-1})$	$(1 + \epsilon^2)\text{opt}_2 + \frac{\epsilon}{k}\text{opt}_\star$	$\text{opt}_F + \epsilon\text{opt}_\star$	$(1 + \epsilon^2)\text{opt}_\star$

Table 1: Comparison of our bounds on the approximation errors of several types of SPSD sketches with those provided in prior works. Only the asymptotically largest terms (as  $\epsilon \rightarrow 0$ ) are displayed and constants are omitted, for simplicity. Here,  $\epsilon \in (0, 1)$ ,  $\text{opt}_\xi$  is the smallest  $\xi$ -norm error possible when approximating  $\mathbf{A}$  with a rank- $k$  matrix ( $k \geq \ln n$ ),  $r = \text{rank}(\mathbf{A})$ ,  $\ell$  is the number of column samples sufficient for the stated bounds to hold,  $k$  is a target rank, and  $\mu_s$  is the coherence of  $\mathbf{A}$  relative to the best rank- $s$  approximation to  $\mathbf{A}$ . The parameter  $\beta \in (0, 1]$  allows for the possibility of sampling using  $\beta$ -approximate leverage scores (see Section 4.2.1) rather than the exact leverage scores. With the exception of (Drineas and Mahoney, 2005), which samples columns with probability proportional to their Euclidean norms, and our novel leverage-based Nyström bound, these bounds are for sampling columns or linear combinations of columns uniformly at random. All bounds hold with constant probability.

practice. In Section 3, we show this second scheme is an instantiation of the power method (as described in Section 2.3) with  $q = 1$ . Accordingly, the deterministic and stochastic error bounds provided in Section 4 provide theoretical guarantees for this SPSD sketch.

Enron, $k = 60$									
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _2 / \ \mathbf{A} - \mathbf{A}_k\ _2$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.386/1.386/1.386	1.386/1.386/1.386	1.386/1.386/1.386	1.386/1.386/1.386		1.570/2.104/2.197	1.496/2.100/2.196	1.023/1.350/2.050	
SRP <sub>T</sub> sketch	1.378/1.379/1.381	1.357/1.360/1.364	1.357/1.360/1.364	1.310/1.317/1.323		1.835/1.950/2.039	1.686/1.874/2.009	1.187/1.287/1.405	
Gaussian sketch	1.378/1.380/1.381	1.357/1.360/1.364	1.357/1.360/1.364	1.314/1.318/1.323		1.812/1.956/2.058	1.653/1.894/2.007	1.187/1.293/1.438	
Leverage sketch	1.321/1.381/1.386	1.039/1.188/1.386	1.039/1.188/1.386	1.039/1.042/1.113		1.345/1.644/2.166	1.198/1.498/2.160	0.942/0.994/1.073	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _F / \ \mathbf{A} - \mathbf{A}_k\ _F$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.004/1.004/1.004	0.993/0.994/0.994	0.993/0.994/0.994	0.972/0.972/0.973		1.041/1.054/1.065	1.023/1.042/1.054	0.867/0.877/0.894	
SRP <sub>T</sub> sketch	1.004/1.004/1.004	0.994/0.994/0.994	0.994/0.994/0.994	0.972/0.972/0.972		1.049/1.054/1.058	1.032/1.037/1.043	0.873/0.877/0.880	
Gaussian sketch	1.004/1.004/1.004	0.994/0.994/0.994	0.994/0.994/0.994	0.972/0.972/0.972		1.049/1.054/1.060	1.032/1.039/1.043	0.874/0.878/0.883	
Leverage sketch	1.002/1.002/1.003	0.994/0.995/1.003	0.994/0.995/1.003	0.988/0.989/0.989		1.027/1.036/1.054	1.011/1.018/1.034	0.862/0.868/0.875	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _* / \ \mathbf{A} - \mathbf{A}_k\ _*$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.002/1.002/1.003	0.984/0.984/0.984	0.984/0.984/0.984	0.943/0.944/0.944		1.011/1.014/1.018	0.988/0.994/0.998	0.760/0.764/0.770	
SRP <sub>T</sub> sketch	1.002/1.002/1.002	0.984/0.984/0.984	0.984/0.984/0.984	0.944/0.944/0.944		1.013/1.015/1.016	0.990/0.993/0.995	0.762/0.764/0.766	
Gaussian sketch	1.002/1.002/1.002	0.984/0.984/0.984	0.984/0.984/0.984	0.944/0.944/0.944		1.013/1.015/1.017	0.991/0.993/0.994	0.762/0.765/0.767	
Leverage sketch	1.002/1.002/1.003	0.990/0.991/0.992	0.990/0.991/0.992	0.977/0.978/0.980		1.004/1.008/1.014	0.982/0.985/0.991	0.758/0.765/0.771	
Abalone <sub>D</sub> , $\sigma = 1.5$ , $k = 20$									
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _2 / \ \mathbf{A} - \mathbf{A}_k\ _2$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	2.168/2.455/2.569	2.022/2.381/2.569	2.022/2.381/2.569	1.823/2.204/2.567		1.989/2.001/2.002	1.987/1.998/2.002	1.739/1.978/2.002	
SRP <sub>T</sub> sketch	2.329/2.416/2.489	2.146/2.249/2.338	2.146/2.249/2.338	1.741/1.840/1.918		1.910/1.938/1.966	1.840/1.873/1.905	1.624/1.669/1.709	
Gaussian sketch	2.347/2.409/2.484	2.161/2.254/2.361	2.161/2.254/2.361	1.723/1.822/1.951		1.903/1.942/1.966	1.839/1.873/1.910	1.619/1.670/1.707	
Leverage sketch	1.508/1.859/2.377	1.152/1.417/2.036	1.152/1.417/2.036	0.774/0.908/1.091		1.242/1.762/1.995	1.000/1.317/1.987	1.000/1.000/1.005	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _F / \ \mathbf{A} - \mathbf{A}_k\ _F$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.078/1.090/1.098	1.061/1.078/1.091	1.061/1.078/1.091	1.026/1.040/1.054		1.036/1.040/1.043	1.028/1.034/1.038	0.998/1.009/1.018	
SRP <sub>T</sub> sketch	1.088/1.089/1.090	1.074/1.075/1.077	1.074/1.075/1.077	1.034/1.035/1.037		1.038/1.039/1.039	1.029/1.030/1.030	1.000/1.000/1.001	
Gaussian sketch	1.087/1.089/1.091	1.073/1.075/1.077	1.073/1.075/1.077	1.033/1.035/1.036		1.038/1.039/1.039	1.029/1.030/1.030	1.000/1.000/1.001	
Leverage sketch	1.028/1.040/1.059	0.998/1.006/1.020	0.998/1.006/1.020	0.959/0.963/0.968		1.004/1.011/1.018	0.996/1.000/1.005	0.994/0.995/0.997	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _* / \ \mathbf{A} - \mathbf{A}_k\ _*$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.022/1.024/1.026	1.010/1.014/1.016	1.010/1.014/1.016	0.977/0.980/0.983		1.013/1.015/1.016	1.002/1.005/1.007	0.965/0.970/0.976	
SRP <sub>T</sub> sketch	1.024/1.024/1.024	1.014/1.014/1.014	1.014/1.014/1.014	0.980/0.980/0.981		1.014/1.014/1.015	1.004/1.004/1.004	0.970/0.970/0.970	
Gaussian sketch	1.024/1.024/1.024	1.014/1.014/1.014	1.014/1.014/1.014	0.980/0.980/0.981		1.014/1.014/1.015	1.004/1.004/1.004	0.970/0.970/0.970	
Leverage sketch	1.009/1.012/1.016	0.994/0.997/1.000	0.994/0.997/1.000	0.965/0.968/0.971		1.002/1.005/1.009	0.997/0.999/1.002	0.995/0.996/0.997	
Protein, $k = 10$									
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _2 / \ \mathbf{A} - \mathbf{A}_k\ _2$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.570/2.104/2.197	1.496/2.100/2.196	1.496/2.100/2.196	1.023/1.350/2.050		1.041/1.054/1.065	1.023/1.042/1.054	0.867/0.877/0.894	
SRP <sub>T</sub> sketch	1.835/1.950/2.039	1.686/1.874/2.009	1.686/1.874/2.009	1.187/1.287/1.405		1.049/1.054/1.058	1.032/1.037/1.043	0.873/0.877/0.880	
Gaussian sketch	1.812/1.956/2.058	1.653/1.894/2.007	1.653/1.894/2.007	1.187/1.293/1.438		1.049/1.054/1.060	1.032/1.039/1.043	0.874/0.878/0.883	
Leverage sketch	1.345/1.644/2.166	1.198/1.498/2.160	1.198/1.498/2.160	0.942/0.994/1.073		1.027/1.036/1.054	1.011/1.018/1.034	0.862/0.868/0.875	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _F / \ \mathbf{A} - \mathbf{A}_k\ _F$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.041/1.054/1.065	1.023/1.042/1.054	1.023/1.042/1.054	0.867/0.877/0.894		1.041/1.054/1.065	1.023/1.042/1.054	0.867/0.877/0.894	
SRP <sub>T</sub> sketch	1.049/1.054/1.058	1.032/1.037/1.043	1.032/1.037/1.043	0.873/0.877/0.880		1.049/1.054/1.058	1.032/1.037/1.043	0.873/0.877/0.880	
Gaussian sketch	1.049/1.054/1.060	1.032/1.039/1.043	1.032/1.039/1.043	0.874/0.878/0.883		1.049/1.054/1.060	1.032/1.039/1.043	0.874/0.878/0.883	
Leverage sketch	1.027/1.036/1.054	1.011/1.018/1.034	1.011/1.018/1.034	0.862/0.868/0.875		1.027/1.036/1.054	1.011/1.018/1.034	0.862/0.868/0.875	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _* / \ \mathbf{A} - \mathbf{A}_k\ _*$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.011/1.014/1.018	0.988/0.994/0.998	0.988/0.994/0.998	0.760/0.764/0.770		1.011/1.014/1.018	0.988/0.994/0.998	0.760/0.764/0.770	
SRP <sub>T</sub> sketch	1.013/1.015/1.016	0.990/0.993/0.995	0.990/0.993/0.995	0.762/0.764/0.766		1.013/1.015/1.016	0.990/0.993/0.995	0.762/0.764/0.766	
Gaussian sketch	1.013/1.015/1.017	0.991/0.993/0.994	0.991/0.993/0.994	0.762/0.765/0.767		1.013/1.015/1.017	0.991/0.993/0.994	0.762/0.765/0.767	
Leverage sketch	1.004/1.008/1.014	0.982/0.985/0.991	0.982/0.985/0.991	0.758/0.765/0.771		1.004/1.008/1.014	0.982/0.985/0.991	0.758/0.765/0.771	
Wine <sub>S</sub> , $\sigma = 1$ , $k = 20$									
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _2 / \ \mathbf{A} - \mathbf{A}_k\ _2$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.989/2.001/2.002	1.987/1.998/2.002	1.987/1.998/2.002	1.739/1.978/2.002		1.989/2.001/2.002	1.987/1.998/2.002	1.739/1.978/2.002	
SRP <sub>T</sub> sketch	1.910/1.938/1.966	1.840/1.873/1.905	1.840/1.873/1.905	1.619/1.670/1.707		1.910/1.938/1.966	1.839/1.873/1.910	1.619/1.670/1.707	
Gaussian sketch	1.903/1.942/1.966	1.839/1.873/1.910	1.839/1.873/1.910	1.619/1.670/1.707		1.903/1.942/1.966	1.839/1.873/1.910	1.619/1.670/1.707	
Leverage sketch	1.242/1.762/1.995	1.000/1.317/1.987	1.000/1.317/1.987	1.000/1.000/1.005		1.242/1.762/1.995	1.000/1.317/1.987	1.000/1.000/1.005	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _F / \ \mathbf{A} - \mathbf{A}_k\ _F$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.036/1.040/1.043	1.028/1.034/1.038	1.028/1.034/1.038	0.998/1.009/1.018		1.036/1.040/1.043	1.028/1.034/1.038	0.998/1.009/1.018	
SRP <sub>T</sub> sketch	1.038/1.039/1.039	1.029/1.030/1.030	1.029/1.030/1.030	1.000/1.000/1.001		1.038/1.039/1.039	1.029/1.030/1.030	1.000/1.000/1.001	
Gaussian sketch	1.038/1.039/1.039	1.029/1.030/1.030	1.029/1.030/1.030	1.000/1.000/1.001		1.038/1.039/1.039	1.029/1.030/1.030	1.000/1.000/1.001	
Leverage sketch	1.004/1.011/1.018	0.996/1.000/1.005	0.996/1.000/1.005	0.994/0.995/0.997		1.004/1.011/1.018	0.996/1.000/1.005	0.994/0.995/0.997	
$\ \mathbf{A} - \mathbf{CW}^T \mathbf{C}^T\ _* / \ \mathbf{A} - \mathbf{A}_k\ _*$									
	$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$		$\ell = k + 8$	$\ell = k \ln k$	$\ell = k \ln k$	$\ell = k \ln n$
Nystrom	1.013/1.015/1.016	1.002/1.005/1.007	1.002/1.005/1.007	0.965/0.970/0.976		1.013/1.015/1.016	1.002/1.005/1.007	0.965/0.970/0.976	
SRP <sub>T</sub> sketch	1.014/1.014/1.015	1.004/1.004/1.004	1.004/1.004/1.004	0.970/0.970/0.970		1.014/1.014/1.015	1.004/1.004/1.004	0.970/0.970/0.970	
Gaussian sketch	1.014/1.014/1.015	1.004/1.004/1.004	1.004/1.004/1.004	0.970/0.970/0.970		1.014/1.014/1.015	1.004/1.004/1.004	0.970/0.970/0.970	
Leverage sketch	1.002/1.005/1.009	0.997/0.999/1.002	0.997/0.999/1.002	0.995/0.996/0.997		1.002/1.005/1.009	0.997/0.999/1.002	0.995/0.996/0.997	

Table 2: The min/mean/max ratios of the errors of several non-rank- $k$  approximations for several of the matrices considered in Table 4. Here  $k$  is the target rank and  $\ell$  is the number of column samples used to form the SPSPD sketches. The min/mean/max ratios were computed using 30 trials for each combination of  $\ell$  and sketching method.

REVISITING THE NYSTRÖM METHOD

source, sketch	pred./obs. spectral error	pred./obs. Frobenius error	pred./obs. trace error
Enron, $k = 60$			
Drineas and Mahoney (2005) nonuniform column sampling	3041.0	66.2	–
Belabbas and Wolfe (2009b) uniform column sampling	–	–	2.0
Kumar et al. (2012) uniform column sampling	331.2	77.7	–
Lemma 5 leverage-based	1287.0	20.5	1.2
Lemma 6 Fourier-based	102.1	42.0	1.6
Lemma 7 Gaussian-based	20.1	7.6	1.4
Lemma 8 uniform column sampling	9.4	285.1	9.5
Protein, $k = 10$			
Drineas and Mahoney (2005), nonuniform column sampling	125.2	18.6	–
Belabbas and Wolfe (2009b), uniform column sampling	–	–	3.6
Kumar et al. (2012), uniform column sampling	35.1	20.5	–
Lemma 5, leverage-based	42.4	6.2	2.0
Lemma 6, Fourier-based	155.0	20.4	3.1
Lemma 7, Gaussian-based	5.7	5.6	2.2
Lemma 8, uniform column sampling	90.0	63.4	14.3
AbaloneD, $\sigma = .15, k = 20$			
Drineas and Mahoney (2005), nonuniform column sampling	360.8	42.5	–
Belabbas and Wolfe (2009b), uniform column sampling	–	–	2.0
Kumar et al. (2012), uniform column sampling	62.0	45.7	–
Lemma 5, leverage-based	235.4	14.1	1.3
Lemma 6, Fourier-based	70.1	36.0	1.7
Lemma 7, Gaussian-based	8.7	8.3	1.3
Lemma 8, uniform column sampling	13.2	166.2	9.0
WineS, $\sigma = 1, k = 20$			
Drineas and Mahoney (2005), nonuniform column sampling	408.4	41.1	–
Belabbas and Wolfe (2009b), uniform column sampling	–	–	2.1
Kumar et al. (2012), uniform column sampling	70.3	44.3	–
Lemma 5, leverage-based	244.6	12.9	1.2
Lemma 6, Fourier-based	94.8	36.0	1.7
Lemma 7, Gaussian-based	11.4	8.1	1.4
Lemma 8, uniform column sampling	13.2	162.2	9.1

Table 3: Comparison of the empirically observed approximation errors to the guarantees provided in this and other works, for several data sets. Each approximation was formed using  $\ell = 6k \ln k$  samples. To evaluate the error guarantees,  $\delta = 1/2$  was taken and all constants present in the statements of the bounds were replaced with ones. The observed errors were taken to be the average errors over 30 runs of the approximation algorithms. The data sets, described in Section 3.1, are representative of several classes of matrices prevalent in machine learning applications.

## 2.5 An Overview of Our Bounds

Our bounds in Table 1 (established as Lemmas 5–8 in Section 4.2) exhibit a common structure: for the spectral and Frobenius norms, we see that the additional error is on a larger scale than the optimal error, and the trace norm bounds all guarantee relative error approximations. This follows from the fact, as detailed in Section 4.1, that low-rank approximations that conform to the SPSD sketching model can be understood as forming column-sample/projection-based approximations to the *square root* of  $\mathbf{A}$ , and thus squaring this approximation yields the resulting approximation to  $\mathbf{A}$ . The squaring process unavoidably results in potentially large additional errors in the case of the spectral and Frobenius norms—whether or not the additional errors are large in practice depends upon the properties of the matrix and the form of stochasticity used in the sampling process. For instance, from our bounds it is clear that Gaussian-based SPSD sketches are expected to have lower additional error in the spectral norm than any of the other sketches considered.

From Table 1, we also see, in the case of uniform Nyström extensions, a necessary dependence on the coherence of the input matrix since columns are sampled uniformly at random. However, we also see that the scales of the additional error of the Frobenius and trace norm bounds are substantially improved over those in prior results. The large additional error in the spectral norm error bound is necessary in the worse case (Gittens, 2012). Lemmas 5, 6 and 7 in Section 4.2—which respectively address leverage-based, Fourier-based, and Gaussian-based SPSD sketches—show that spectral norm additive-error bounds with additional error on a substantially smaller scale can be obtained if one first mixes the columns before sampling from  $\mathbf{A}$  or one samples from a judicious nonuniform distribution over the columns.

Table 2 compares the minimum, mean, and maximum approximation errors of several SPSD sketches of four matrices (described in Section 3.1) to the optimal rank- $k$  approximation errors. We consider three regimes for  $\ell$ , the number of column samples used to construct the sketch:  $\ell = O(k)$ ,  $\ell = O(k \ln k)$ , and  $\ell = O(k \ln n)$ . These matrices exhibit a diverse range of properties: e.g., Enron is sparse and has a slowly decaying spectrum, while Protein is dense and has a rapidly decaying spectrum. Yet we notice that the sketches perform quite well on each of these matrices. In particular, when  $\ell = O(k \ln n)$ , the average errors of the sketches are within  $1 + \epsilon$  of the optimal rank- $k$  approximation errors, where  $\epsilon \in [0, 1]$ . Also note that the leverage-based sketches consistently have lower average errors (in all of the three norms considered) than all other sketches. Likewise, the uniform Nyström extensions usually have larger average errors than the other sketches. These two sketches represent opposite extremes: uniform Nyström extensions (constructed using uniform column sampling) are constructed using no knowledge about the matrix, while leverage-based sketches use an importance sampling distribution derived from the SVD of the matrix to determine which columns to use in the construction of the sketch.

Table 3 illustrates the gap between the theoretical results currently available in the literature and what is observed in practice: it depicts the ratio between the error bounds in Table 1 and the average errors observed over 30 runs of the SPSD approximation algorithms (the error bound from (Talwalkar and Rostamizadeh, 2010) is not considered in the table, as it does not apply at the number of samples  $\ell$  used in the experiments). Several trends can be identified; among them, we note that the bounds provided in this paper for Gaussian-based

sketches come quite close to capturing the errors seen in practice, and the Frobenius and trace norm error guarantees of the leverage-based and Fourier-based sketches tend to more closely reflect the empirical behavior than the error guarantees provided in prior work for Nyström sketches. Overall, the trace norm error bounds are quite accurate. On the other hand, prior bounds are sometimes more informative in the case of the spectral norm (with the notable exception of the Gaussian sketches). Several important points can be gleaned from these observations. First, the accuracy of the Gaussian error bounds suggests that the main theoretical contribution of this work, the deterministic structural results given as Theorems 2 through 4, captures the underlying behavior of the SPSD sketching process. This supports our belief that this work provides a foundation for truly informative error bounds. Given that this is the case, it is clear that the analysis of the stochastic elements of the SPSD sketching process is much sharper in the Gaussian case than in the leverage-score, Fourier, and uniform Nyström cases. We expect that, at least in the case of leverage and Fourier-based sketches, the stochastic analysis can and will be sharpened to produce error guarantees almost as informative as the ones we have provided for Gaussian-based sketches.

### 3. Empirical Aspects of SPSD Low-rank Approximation

In this section, we present our main empirical results, which consist of evaluating sampling and projection algorithms applied to a diverse set of SPSD matrices. The bulk of our empirical evaluation considers two random projection procedures and two random sampling procedures for the sketching matrix  $\mathbf{S}$ : for random projections, we consider using SRFTs (Subsampled Randomized Fourier Transforms) as well as uniformly sampling from Gaussian mixtures of the columns; and for random sampling, we consider sampling columns uniformly at random as well as sampling columns according to a nonuniform importance sampling distribution that depends on the empirical statistical leverage scores. In the latter case of leverage score-based sampling, we also consider the use of both the (naïve and expensive) exact algorithm as well as a (recently-developed fast) approximation algorithm. Section 3.1 starts with a brief description of the data sets we consider; Section 3.2 describes the details of our SPSD sketching algorithms; Section 3.3 summarizes our experimental results to help guide in the selection of sketching methods; in Section 3.4, we present our main results on reconstruction quality for the random sampling and random projection methods; and, in Section 3.5, we discuss running time issues, and we present our main results for running time and reconstruction quality for both exact and approximate versions of leverage-based sampling.

We emphasize that we don’t intend these results to be “comprehensive” but instead to be “illustrative” case-studies—that are representative of a much wider range of applications than have been considered previously. In particular, we would like to illustrate the tradeoffs between these methods in different realistic applications in order, *e.g.*, to provide directions for future work. In addition to clarifying some of these issues, our empirical evaluation also illustrates ways in which existing theory is insufficient to explain the success of sampling and projection methods. This motivates our improvements to existing theory that we describe in Section 4.

All of our computations were conducted using 64-bit MATLAB R2012a under Ubuntu on a 2.6-GHz quad-core Intel i7 machine with 6Gb of RAM. To allow for accurate timing

comparisons, all computations were carried out in a single thread. When applied to an  $n \times n$  SPSD matrix  $\mathbf{A}$ , our implementation of the SRFT requires  $O(n^2 \ln n)$  operations, as it applies MATLAB’s `fft` to the entire matrix  $\mathbf{A}$  and *then* it samples  $\ell$  columns from the resulting matrix. A more rigorous implementation of the SRFT algorithm could reduce this running time to  $O(n^2 \ln \ell)$ , but due to the complexities involved in optimizing pruned FFT codes, we did not pursue this avenue.

### 3.1 Data Sets

Table 4 provides summary statistics for the data sets used in our empirical evaluation. We consider four classes of matrices commonly encountered in machine learning and data analysis applications: normalized Laplacians of very sparse graphs drawn from “informatics graph” applications; dense matrices corresponding to Linear Kernels from machine learning applications; dense matrices constructed from a Gaussian Radial Basis Function Kernel (RBFK); and sparse RBFK matrices constructed using Gaussian radial basis functions, truncated to be nonzero only for nearest neighbors. This collection of data sets represents a wide range of data sets with very different (sparsity, spectral, leverage score, etc.) properties that have been of interest recently not only in machine learning but in data analysis more generally.

To understand better the Laplacian data, recall that, given an undirected graph with weighted adjacency matrix  $\mathbf{W}$ , its normalized graph Laplacian is

$$\mathbf{A} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2},$$

where  $\mathbf{D}$  is the diagonal matrix of weighted degrees of the nodes of the graph, *i.e.*,  $D_{ii} = \sum_{j \neq i} W_{ij}$ .

The remaining data sets are kernel matrices associated with data drawn from a variety of application areas. Recall that, given given points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  and a function  $\kappa : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , the  $n \times n$  matrix with elements

$$A_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

is called the kernel matrix of  $\kappa$  with respect to  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . Appropriate choices of  $\kappa$  ensure that  $\mathbf{A}$  is positive semidefinite. When this is the case, the entries  $A_{ij}$  can be interpreted as measuring, in a sense determined by the choice of  $\kappa$ , the similarity of points  $i$  and  $j$ . Specifically, if  $\mathbf{A}$  is SPSD, then  $\kappa$  determines a so-called *feature map*  $\Phi_\kappa : \mathbb{R}^d \rightarrow \mathbb{R}^n$  such that

$$A_{ij} = \langle \Phi_\kappa(\mathbf{x}_i), \Phi_\kappa(\mathbf{x}_j) \rangle$$

measures the similarity (correlation) of  $\mathbf{x}_i$  and  $\mathbf{x}_j$  in feature space (Schölkopf and Smola, 2001).

When  $\kappa$  is the usual Euclidean inner-product, so that

$$A_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle,$$

$\mathbf{A}$  is called a Linear Kernel matrix. Gaussian RBFK matrices, defined by

$$A_{ij}^\sigma = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}\right),$$

Name	Description	n	d	%nnz
Laplacian Kernels				
HEP	arXiv High Energy Physics collaboration graph	9877	NA	0.06
GR	arXiv General Relativity collaboration graph	5242	NA	0.12
Enron	subgraph of the Enron email graph	10000	NA	0.22
Gnutella	Gnutella peer to peer network on Aug. 6, 2002	8717	NA	0.09
Linear Kernels				
Dexter	bag of words	2000	20000	83.8
Protein	derived feature matrix for <i>S. cerevisiae</i>	6621	357	99.7
SNPs	DNA microarray data from cancer patients	5520	43	100
Gisette	images of handwritten digits	6000	5000	100
Dense RBF Kernels				
AbaloneD	physical measurements of abalones	4177	8	100
WineD	chemical measurements of wine	4898	12	100
Sparse RBF Kernels				
AbaloneS	physical measurements of abalones	4177	8	82.9/48.1
WineS	chemical measurements of wine	4898	12	11.1/88.0

Table 4: The data sets used in our empirical evaluation (Leskovec et al., 2007; Klimt and Yang, 2004; Guyon et al., 2005; Gustafson et al., 2006; Nielsen et al., 2002; Corke, 1996; Asuncion and Newman, 2012). Here,  $n$  is the number of data points,  $d$  is the number of features in the input space before kernelization, and %nnz is the percentage of nonzero entries in the matrix. For Laplacian “kernels,”  $n$  is the number of nodes in the graph (and thus there is no  $d$  since the graph is “given” rather than “constructed”). The %nnz for the Sparse RBF Kernels depends on the  $\sigma$  parameter; see Table 5.

correspond to the similarity measure  $\kappa(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|_2^2/\sigma^2)$ . Here  $\sigma$ , a nonnegative number, defines the scale of the kernel. Informally,  $\sigma$  defines the “size scale” over which pairs of points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  “see” each other. Typically  $\sigma$  is determined by a global cross-validation criterion, as  $\mathbf{A}^\sigma$  is generated for some specific machine learning task; and, thus, one may have no *a priori* knowledge of the behavior of the spectrum or leverage scores of  $\mathbf{A}^\sigma$  as  $\sigma$  is varied. Accordingly, we consider Gaussian RBFK matrices with different values of  $\sigma$ .

Finally, given the same data points,  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , one can construct sparse Gaussian RBFK matrices

$$A_{ij}^{(\sigma, \nu, C)} = \left[ \left( 1 - \frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2}{C} \right)^\nu \right]^+ \cdot \exp \left( -\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2} \right),$$

where  $[x]^+ = \max\{0, x\}$ . When  $\nu$  is larger than  $(d + 1)/2$ , this kernel matrix is positive semidefinite (Genton, 2002). Increasing  $\nu$  shrinks the magnitudes of the off-diagonal entries of the matrix toward zero. As the cutoff point  $C$  decreases the matrix becomes more sparse; in particular,  $C \rightarrow 0$  ensures that  $\mathbf{A}^{(\sigma, \nu, C)} \rightarrow \mathbf{I}$ . On the other hand,  $C \rightarrow \infty$  ensures that

Name	%nnz	$\lceil \frac{\ \mathbf{A}\ _F^2}{\ \mathbf{A}\ _2^2} \rceil$	$k$	$\frac{\lambda_{k+1}}{\lambda_k}$	$100 \frac{\ \mathbf{A} - \mathbf{A}_k\ _F}{\ \mathbf{A}\ _F}$	$100 \frac{\ \mathbf{A} - \mathbf{A}_k\ _*}{\ \mathbf{A}\ _*}$	$k$ th-largest leverage score scaled by $n/k$
HEP	0.06	3078	20	0.998	7.8	0.4	128.8
HEP	0.06	3078	60	0.998	13.2	1.1	41.9
GR	0.12	1679	20	0.999	10.5	0.74	71.6
GR	0.12	1679	60	1	17.9	2.16	25.3
Enron	0.22	2588	20	0.997	7.77	0.352	245.8
Enron	0.22	2588	60	0.999	12.0	0.94	49.6
Gnutella	0.09	2757	20	1	8.1	0.41	166.2
Gnutella	0.09	2757	60	0.999	13.7	1.20	49.4
Dexter	83.8	176	8	0.963	14.5	.934	16.6
Protein	99.7	24	10	0.987	42.6	7.66	5.45
SNPs	100	3	5	0.928	85.5	37.6	2.64
Gisette	100	4	12	0.90	90.1	14.6	2.46
AbaloneD (dense, $\sigma = .15$ )	100	41	20	0.992	42.1	3.21	18.11
AbaloneD (dense, $\sigma = 1$ )	100	4	20	0.935	97.8	59	2.44
WineD (dense, $\sigma = 1$ )	100	31	20	0.99	43.1	3.89	26.2
WineD (dense, $\sigma = 2.1$ )	100	3	20	0.936	94.8	31.2	2.29
AbaloneS (sparse, $\sigma = .15$ )	82.9	400	20	0.989	15.4	1.06	48.4
AbaloneS (sparse, $\sigma = 1$ )	48.1	5	20	0.982	90.6	21.8	3.57
WineS (sparse, $\sigma = 1$ )	11.1	116	20	0.995	29.5	2.29	49.0
WineS (sparse, $\sigma = 2.1$ )	88.0	39	20	0.992	41.6	3.53	24.1

Table 5: Summary statistics for the data sets from Table 4 that we used in our empirical evaluation.

$\mathbf{A}^{(\sigma, \nu, C)}$  approaches the (dense) Gaussian RBFK matrix  $\mathbf{A}^\sigma$ . For simplicity, in our empirical evaluations, we fix  $\nu = \lceil (d+1)/2 \rceil$  and  $C = 3\sigma$ , and we vary  $\sigma$ .

To illustrate the diverse range of properties exhibited by these four classes of data sets, consider Table 5. Several observations are particularly relevant to our discussion below.

- All of the Laplacian Kernels drawn from informatics graph applications are extremely sparse in terms of number of nonzeros, and they all tend to have very slow spectral decay, as illustrated both by the quantity  $\lceil \|\mathbf{A}\|_F^2 / \|\mathbf{A}\|_2^2 \rceil$  (this is the *stable rank*, which is a numerically stable (under)estimate of the rank of  $\mathbf{A}$ ) as well as by the relatively small fraction of the Frobenius norm that is captured by the best rank- $k$  approximation to  $\mathbf{A}$ .
- Both the Linear Kernels and the Dense RBF Kernels are much denser and are much more well-approximated by moderately to very low-rank matrices. In addition, both the Linear Kernels and the Dense RBF Kernels have statistical leverage scores that are much more uniform—there are several ways to illustrate this, none of them perfect. Here, we illustrate this by considering the  $k^{\text{th}}$  largest leverage score, scaled by the factor  $n/k$  (if  $\mathbf{A}$  were exactly rank  $k$ , this would be the coherence of  $\mathbf{A}$ ). For the Linear Kernels and the Dense RBF Kernels, this quantity is typically one to two orders of magnitude smaller than for the Laplacian Kernels.



- For the Dense RBF Kernels, we consider two values of the  $\sigma$  parameter, again chosen (somewhat) arbitrarily. For both AbaloneD and WineD, we see that decreasing  $\sigma$  from 1 to 0.15, *i.e.*, letting data points “see” fewer nearby points, has two important effects: first, it results in matrices that are much *less* well-approximated by low-rank matrices; and second, it results in matrices that have *much* more heterogeneous leverage scores.
- For the Sparse RBF Kernels, there are a range of sparsities, ranging from above the sparsity of the sparsest Linear Kernel, but all are denser than the Laplacian Kernels. Changing the  $\sigma$  parameter has the same effect (although it is even more pronounced) for Sparse RBF Kernels as it has for Dense RBF Kernels. In addition, “sparsifying” a Dense RBF Kernel also has the effect of making the matrix less well approximated by a low-rank matrix and of making the leverage scores more nonuniform.

As we see below, when we consider the RBF Kernels as the width parameter and sparsity are varied, we observe a range of intermediate cases between the extremes of the (“nice”) Linear Kernels and the (very “non-nice”) Laplacian Kernels.

### 3.2 SPSD Sketching Algorithms

The sketching matrix  $\mathbf{S}$  may be selected in a variety of ways. For sampling-based sketches, the sketching matrix  $\mathbf{S}$  contains exactly one nonzero in each column, corresponding to a single sample from the columns of  $\mathbf{A}$ . For projection-based sketches,  $\mathbf{S}$  is dense, and mixes the columns of  $\mathbf{A}$  before sampling from the resulting matrix.

In more detail, we consider two types of sampling-based SPSD sketches (*i.e.* Nyström extensions): those constructed by sampling columns uniformly at random with replacement, and those constructed by sampling columns from a distribution based upon the leverage scores of the matrix filtered through the optimal rank- $k$  approximation of the matrix. In the case of column sampling, the sketching matrix  $\mathbf{S}$  is simply the first  $\ell$  columns of a matrix that was chosen uniformly at random from the set of all permutation matrices.

In the case of leverage-based sampling,  $\mathbf{S}$  has a more complicated distribution. Recall that the leverage scores relative to the best rank- $k$  approximation to  $\mathbf{A}$  are the squared Euclidean norms of the rows of the  $n \times k$  matrix  $\mathbf{U}_1$  :

$$\ell_j = \|(\mathbf{U}_1)_j\|^2.$$

It follows from the orthonormality of  $\mathbf{U}_1$  that  $\sum_j (\ell_j/k) = 1$ , and the leverage scores can thus be interpreted as a probability distribution over the columns of  $\mathbf{A}$ . To construct a sketching matrix corresponding to sampling from this distribution, we first select the columns to be used by sampling with replacement from this distribution. Then,  $\mathbf{S}$  is constructed as  $\mathbf{S} = \mathbf{R}\mathbf{D}$  where  $\mathbf{R} \in \mathbb{R}^{n \times \ell}$  is a column selection matrix that samples columns of  $\mathbf{A}$  from the given distribution—*i.e.*,  $\mathbf{R}_{ij} = 1$  iff the  $i$ th column of  $\mathbf{A}$  is the  $j$ th column selected—and  $\mathbf{D}$  is a diagonal rescaling matrix satisfying  $\mathbf{D}_{jj} = \frac{1}{\sqrt{\ell p_i}}$  iff  $\mathbf{R}_{ij} = 1$ . Here,  $p_i = \ell_i/k$  is the probability of choosing the  $i$ th column of  $\mathbf{A}$ . It is often expensive to compute the leverage scores exactly; in Section 3.5, we consider the performance of sketches based on several leverage score approximation algorithms.

The two projection-based sketches we consider use Gaussians and the real Fourier transform. In the former case,  $\mathbf{S}$  is a matrix of i.i.d.  $\mathcal{N}(0, 1)$  random variables. In the latter case,

$\mathbf{S}$  is a *subsampled randomized Fourier transform* (SRFT) matrix; that is,  $\mathbf{S} = \sqrt{\frac{n}{\ell}} \mathbf{D} \mathbf{F} \mathbf{R}$ , where  $\mathbf{D}$  is a diagonal matrix of Rademacher random variables,  $\mathbf{F}$  is the real Fourier transform matrix, and  $\mathbf{R}$  restricts to  $\ell$  columns.

For conciseness, we do not present results for sampling-based sketches where rows are selected with probability proportional to their row norms. This form of sampling can be similar to leverage-score sampling for sparse graphs with highly connected vertices (Mahoney and Drineas, 2009), and in cases where the matrix has been preprocessed to have uniform row lengths, reduces to uniform sampling.

In the figures, we refer to sketches constructed by selecting columns uniformly at random with the label ‘unif’, leverage score-based sketches with ‘lev’, Gaussian sketches with ‘gaussian’, and Fourier sketches with ‘srft’.

### 3.3 Guidelines for Selecting Sketching Schemes

In the remainder of this section of the paper, we provide empirical evaluations of the sampling and projection-based sketching schemes just described, with an eye towards identifying the aspects of the datasets that affect the relative performance of the sketching schemes. However our experiments also provide some practical guidelines for selecting a particular sketching scheme.

- Despite the theoretical result that the worst-case spectral error in using Nyström sketches obtained via uniform column-samples can be much worse than that of using projection or leverage-based sketches, on the corpus of data sets we considered, such sketches perform within a small multiple of the error of more computationally expensive leverage-based and projection-based sketches. For data sets with more nonuniform leverage score properties, random projections and leverage-based sampling will do better (Ma et al., 2014).
- In the case where parsimony of the sketch is of primary concern, *i.e.* where the primary concern is to maintain  $\ell \approx k$ , leverage sketches are an attractive option. In particular, when an RBF kernel with small bandwidth is used, or the data set is sparse, leverage-based sketches often provide higher accuracy than projection or uniform-sampling based sketches.
- The norm in which the error is measured should be taken into consideration when selecting the sketching algorithm. In particular, sketches which use power iterations are most useful when the error is measured in the spectral norm, and in this case, projection-based sketches (in particular, *prolonged* sketches—see Section 3.6) noticeably outperform uniform sampling-based sketches.

### 3.4 Reconstruction Accuracy of Sampling and Projection Algorithms

Here, we describe the performances of the SPSD sketches described in Section 3.2—column sampling uniformly at random without replacement, column sampling according to the nonuniform leverage score probabilities, and sampling using Gaussian and SRFT mixtures of the columns—in terms of reconstruction accuracy for the data sets described in Section 3.1. We describe general observations we have made about each class of matrices in

turn, and then we summarize our observations. We consider only the use of exact leverage scores here, and we postpone until Section 3.5 a discussion of running time issues and similar reconstruction results when approximate leverage scores are used for the importance sampling distribution. The relative errors

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_\xi / \|\mathbf{A} - \mathbf{A}_k\|_\xi \quad (7)$$

are plotted, with each point in the figures of this section representing the average errors observed over 30 trials.

### 3.4.1 GRAPH LAPLACIANS

Figure 1 and Figure 2 show the reconstruction error results for sampling and projection methods applied to several normalized graph Laplacians. The former shows GR and HEP, each for two values of the rank parameter, and the latter shows Enron and Gnutella, again each for two values of the rank parameter. Both figures show the spectral, Frobenius, and trace norm approximation errors, as a function of the number of column samples  $\ell$ , relative to the error of the optimal rank- $k$  approximation of  $\mathbf{A}$ .

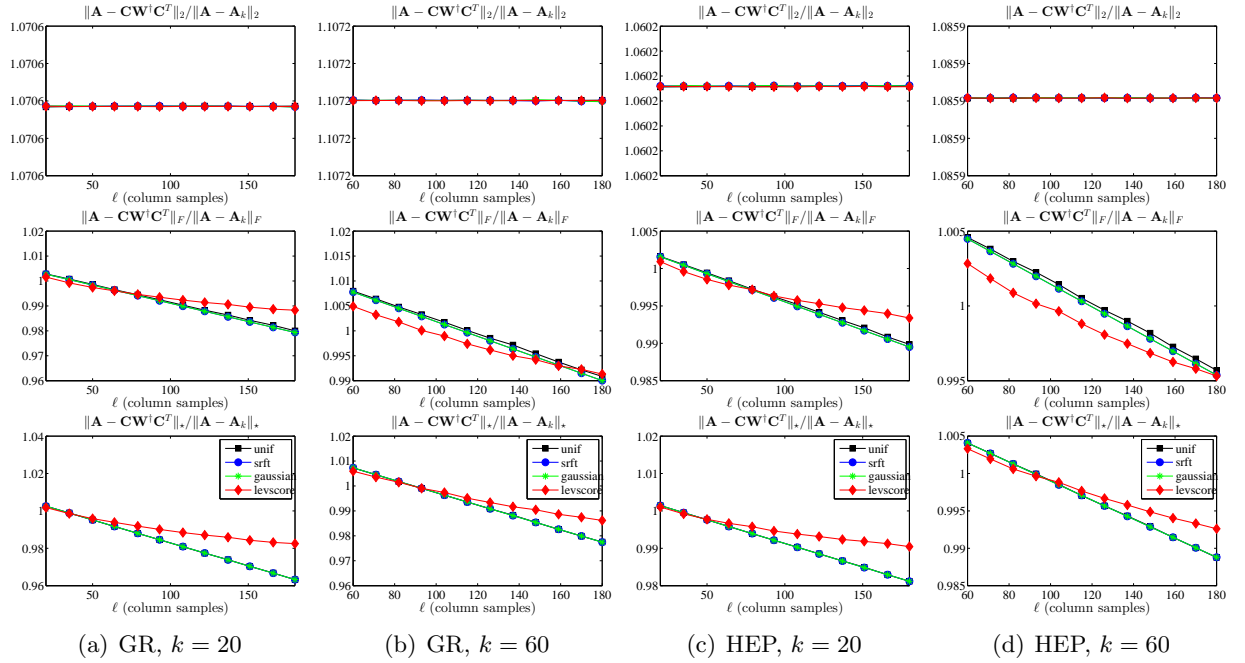


Figure 1: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several SPSD sketches, as a function of the number of column samples  $\ell$ , for the GR and HEP Laplacian data sets, with two choices of the rank parameter  $k$ .

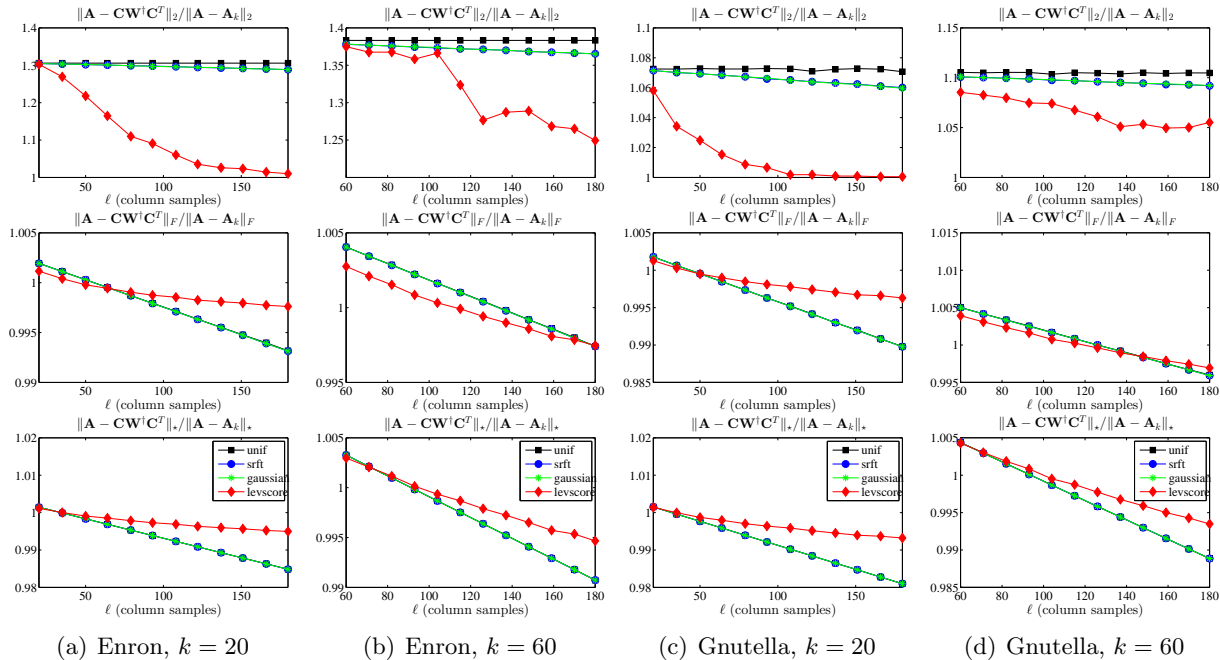


Figure 2: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several SPSP sketches, as a function of the number of column samples  $\ell$ , for the Enron and Gnutella Laplacian data sets, with two choices of the rank parameter  $k$ .

These and subsequent figures contain a lot of information, some of which is peculiar to the given data sets and some of which is more general. In light of subsequent discussion, several observations are worth making about the results presented in these two figures.

- All of the SPSP sketches provide quite accurate approximations—relative to the best possible approximation factor for that norm, and relative to bounds provided by existing theory, as reviewed in Section 2.4—even with only  $k$  column samples (or in the case of the Gaussian and SRFT mixtures, with only  $k$  linear combinations of columns). Upon examination, this is partly due to the extreme sparsity and extremely slow spectral decay of these data sets which means, as shown in Table 4, that only a small fraction of the (spectral or Frobenius or trace) mass is captured by the optimal rank 20 or 60 approximation. Thus, although an SPSP sketch constructed from 20 or 60 vectors also only captures a small portion of the mass of the matrix, the relative error is small, since the scale of the residual error is large.
- The scale of the Y axes is different between different figures and subfigures. This is to highlight properties within a given plot, but it can hide several things. In particular, note that the scale for the spectral norm is generally larger than for the Frobenius norm, which is generally larger than for the trace norm, consistent with the size of those norms; and that the scale is larger for higher-rank approximations, *e.g.* compare

GR  $k = 20$  with GR  $k = 60$ . This is also consistent with the larger amount of mass captured by higher-rank approximations.

- For  $\ell > k$ , the errors tend to decrease (or at least not increase, as for GR and HEP the spectral norm error is flat as a function of  $\ell$ ), which is intuitive.
- The X axes ranges from  $k$  to  $9k$  for the  $k = 20$  plots and from  $k$  to  $3k$  for the  $k = 60$  plots. As a practical matter, choosing  $\ell$  between  $k$  and (say)  $2k$  or  $3k$  is probably of greatest interest. In this regime, there is an interesting tradeoff: for moderately large values of  $\ell$  in this regime, the error for leverage-based sampling is moderately better than for uniform sampling or random projections, while if one chooses  $\ell$  to be much larger then the improvements from leverage-based sampling saturate and the uniform sampling and random projection methods are better. This is most obvious in the Frobenius norm plots, although it is also seen in the trace norm plots, and it suggests that some combination of leverage-based sampling and uniform sampling might be best.
- The behavior of the approximations with respect to the spectral norm is quite different from the behavior in the Frobenius and trace norms. In the latter, as the number of samples  $\ell$  increases, the errors tend to decrease; while for the former, the errors tend to be much flatter as a function of increasing  $\ell$  for at least the Gaussian, SRFT, and uniformly sampled sketches.

All in all, there seems to be quite complicated behavior for low-rank sketches for these Laplacian data sets. Several of these observations can also be made for subsequent figures; but in some other cases the (very sparse and not very low rank) structural properties of the data are primarily responsible.

### 3.4.2 LINEAR KERNELS

Figure 3 shows the reconstruction error results for sampling and projection methods applied to several Linear Kernels. The data sets (Dexter, Protein, SNPs, and Gisette) are all quite low-rank and have fairly uniform leverage scores. Several observations are worth making about the results presented in this figure.

- All of the methods perform quite similarly: all have errors that decrease smoothly with increasing  $\ell$ , and in this case there is little advantage to using methods other than uniform sampling (since they perform similarly and are more expensive). Also, since the ranks are so low and the leverage scores are so uniform, the leverage score sketch is no longer significantly distinguished by its tendency to saturate quickly.
- The scale of the Y axes is much larger than for the Laplacian data sets, mostly since the matrices are much more well-approximated by low-rank matrices, although the scale decreases as one goes from spectral to Frobenius to trace reconstruction error, as before.

These linear kernels (and also to some extent the dense RBF kernels below that have larger  $\sigma$  parameter) are examples of relatively “nice” machine learning data sets that are similar

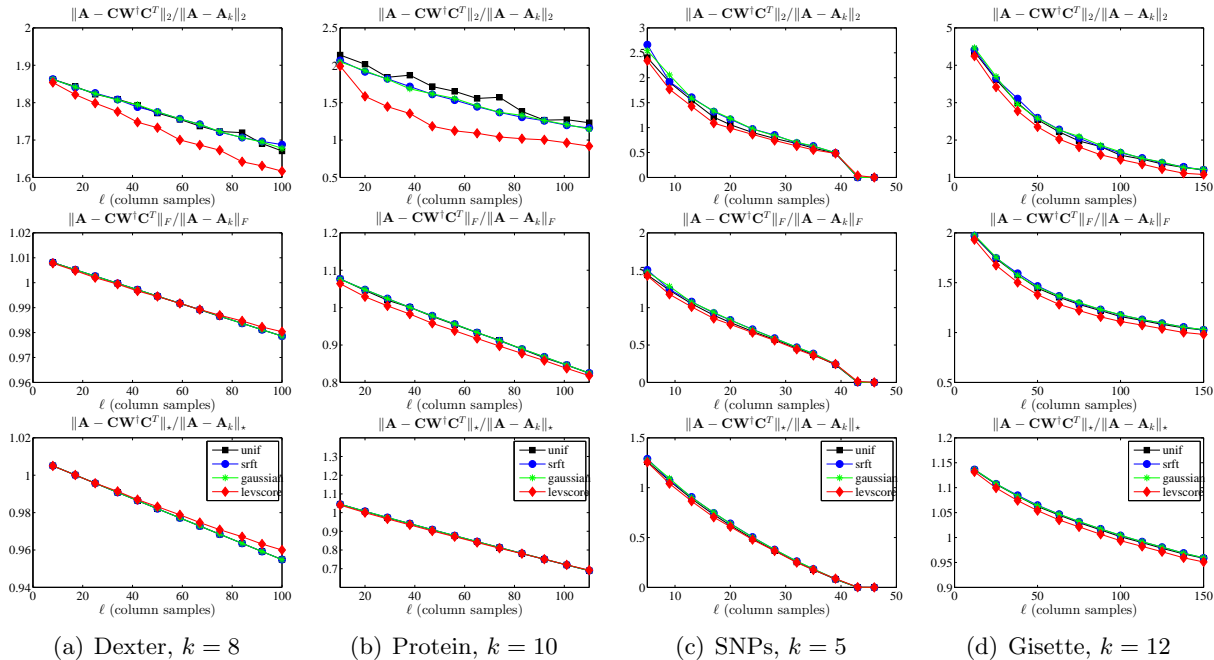


Figure 3: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several SPSP sketches, as a function of the number of column samples  $\ell$ , for the Linear Kernel data sets.

to matrices where uniform sampling has been shown to perform well previously (Talwalkar et al., 2008; Kumar et al., 2009a,c, 2012); for these matrices our empirical results agree with these prior works.

### 3.4.3 DENSE AND SPARSE RBF KERNELS

Figure 4 and Figure 5 present the reconstruction error results for sampling and projection methods applied to several dense RBF and sparse RBF kernels. Several observations are worth making about the results presented in these figures.

- All of the methods have errors that decrease with increasing  $\ell$ , but for larger values of  $\sigma$  and for denser data, the decrease is somewhat more regular, and the four methods tend to perform similarly. For larger values of  $\sigma$  and sparser data, leverage score sampling is somewhat better. This parallels what we observed with the Linear Kernels, except that here the leverage score sampling is somewhat better for all values of  $\ell$ .
- For smaller values of  $\sigma$ , leverage score sampling tends to be much better than uniform sampling and projection-based methods. For sparse data, however, this effect saturates; and we again observe (especially when  $\sigma$  is smaller in AbaloneS and WineS) the tradeoff we observed previously with the Laplacian data—leverage score sampling is better when  $\ell$  is moderately larger than  $k$ , while uniform sampling and random projections are better when  $\ell$  is much larger than  $k$ .

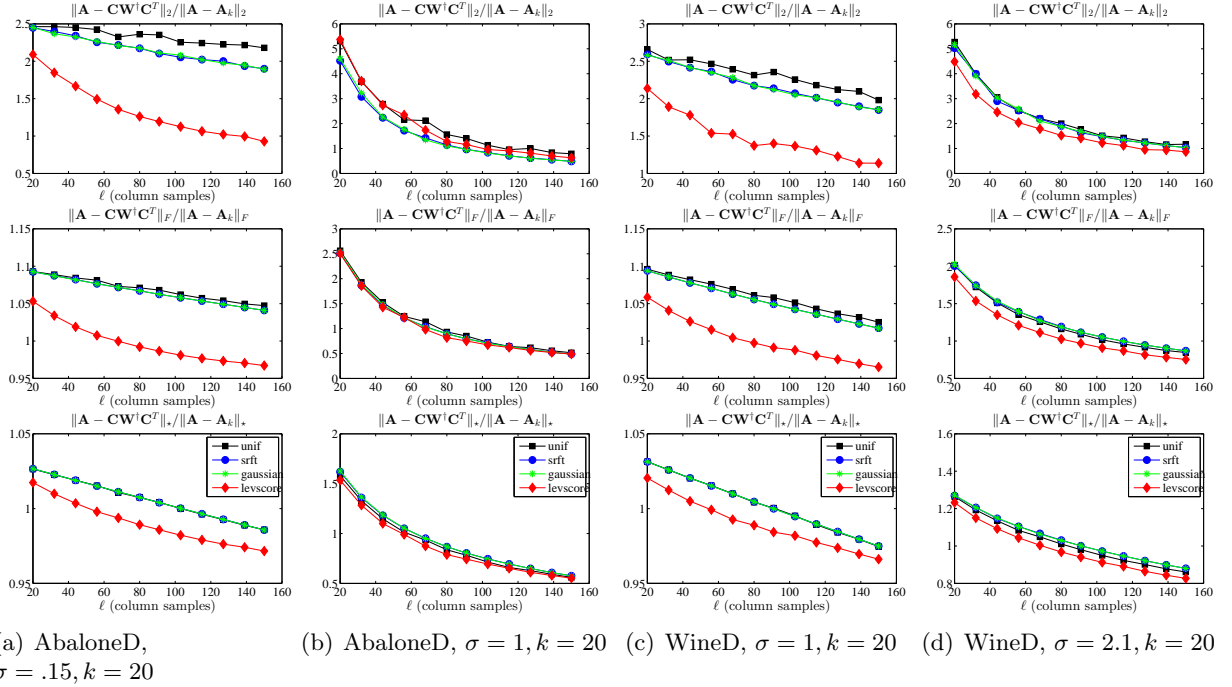


Figure 4: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several SPSD sketches, as a function of the number of column samples  $\ell$ , for several dense RBF data sets.

Recall from Table 5 that for smaller values of  $\sigma$  and for sparser kernels, the SPSD matrices are less well-approximated by low-rank matrices, and they have more heterogeneous leverage scores. Thus, they are more similar to the Laplacian data than the Linear Kernel data; this suggests (as we have observed) that leverage score sampling should perform relatively better than uniform column sampling and projection-based schemes when in these two cases.

#### 3.4.4 SUMMARY OF COMPARISON OF SAMPLING AND PROJECTION ALGORITHMS

Before proceeding, there are several summary observations that we can make about sampling versus projection methods for the data sets we have considered.

- Linear Kernels and to a lesser extent Dense RBF Kernels with larger  $\sigma$  parameter have relatively low rank and relatively uniform leverage scores, and in these cases uniform sampling does quite well. These data sets correspond most closely with those that have been studied previously in the machine learning literature, and for these data sets our results are in agreement with that prior work.
- Sparsifying RBF Kernels and/or choosing a smaller  $\sigma$  parameter tends to make these kernels less well-approximated by low-rank matrices and to have more heterogeneous leverage scores. In general, these two properties need not be directly related—the spectrum is a property of eigenvalues, while the leverage scores are determined by the

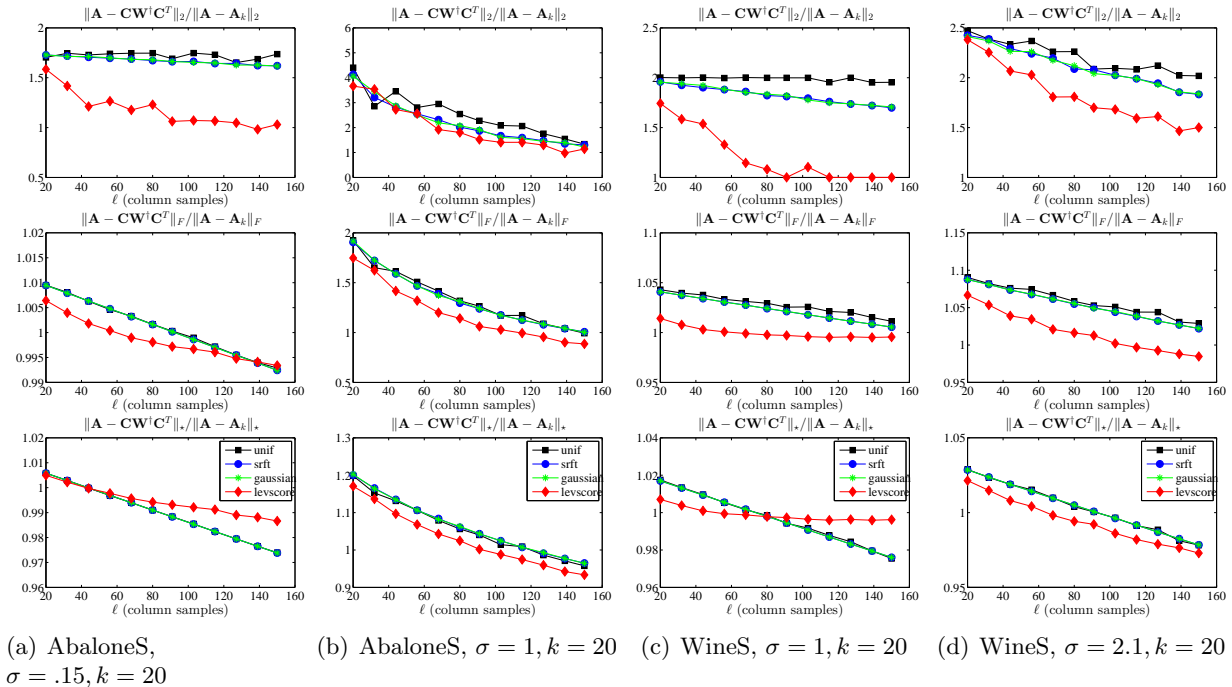


Figure 5: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several SPSD sketches, as a function of the number of column samples  $\ell$ , for several sparse RBF data sets.

eigenvectors—but for the data we examined they are related, in that matrices with more slowly decaying spectra also often have more heterogeneous leverage scores.

- For Dense RBF Kernels with smaller  $\sigma$  and Sparse RBF Kernels, leverage score sampling tends to do much better than other methods. Interestingly, the Sparse RBF Kernels have many properties of very sparse Laplacian Kernels corresponding to relatively-unstructured informatics graphs, an observation which should be of interest for researchers who construct sparse graphs from data using, *e.g.*, “locally linear” methods, to try to reconstruct hypothesized low-dimensional manifolds.
- Reconstruction quality under leverage score sampling saturates, as a function of choosing more samples  $\ell$ . As a consequence, there can be a tradeoff between leverage score sampling or other methods being better, depending on the values of  $\ell$  that are chosen.

In general, *all* of the sampling and projection methods we considered perform *much* better on the SPSD matrices we considered than previous worst-case bounds (*e.g.*, (Drineas and Mahoney, 2005; Kumar et al., 2012; Gittens, 2012)) would suggest. Specifically, even the worst results correspond to single-digit approximation factors in relative scale. This observation is intriguing, because the motivation of leverage score sampling (recall that in this context random projections should be viewed as performing uniform random sampling in a



randomly-rotated basis where the leverage scores have been approximately uniformized (Mahoney, 2011)) is very much tied to the Frobenius norm, and so there is no *a priori* reason to expect its good performance to extend to the spectral or trace norms. Motivated by this, we revisit the question of proving improved worst-case theoretical bounds in Section 4.

Before describing these improved theoretical results, however, we address in Section 3.5 running time questions. After all, a naïve implementation of sampling with exact leverage scores is slower than other methods (and much slower than uniform sampling). As shown below, by using the recently-developed approximation algorithm of Drineas et al. (2012), not only does this approximation algorithm run in time comparable with random projections (for certain parameter settings), it also leads to approximations that soften the strong bias that the exact leverage scores provide toward the best rank- $k$  approximation to the matrix, thereby leading to improved reconstruction results in many cases.

### 3.5 Reconstruction Accuracy of Leverage Score Approximation Algorithms

A naïve view might assume that computing probabilities that permit leverage-based sampling requires an  $O(n^3)$  computation of the full SVD, or at least the full computation of a partial SVD, and thus that it would be much more expensive than recently-developed random projection methods. Indeed, an “exact” computation of the leverage scores with a truncated SVD takes roughly  $O(n^2k)$  time. Recent work, however, has shown that relative-error approximations to all the statistical leverage scores can be computed more quickly than this exact algorithm (Drineas et al., 2012). Here, we implement and evaluate a version of this algorithm. We evaluate it both in terms of running time and in terms of reconstruction quality on the diverse suite of real data matrices we considered above. This is the first work to provide an empirical evaluation of an implementation of the leverage score approximation algorithms of Drineas et al. (2012), illustrating empirically the tradeoffs between cost and efficiency in a practical setting.

#### 3.5.1 DESCRIPTION OF THE FAST APPROXIMATION ALGORITHM OF DRINEAS ET AL. (2012)

Algorithm 1 (which originally appeared as Algorithm 1 in Drineas et al. (2012)) takes as input an arbitrary  $n \times d$  matrix  $\mathbf{A}$ , where  $n \gg d$ , and it returns as output a  $1 \pm \epsilon$  approximation to *all* of the statistical leverage scores of the input matrix. The original algorithm of Drineas et al. (2012) uses a subsampled Hadamard transform and requires  $r_1$  to be somewhat larger than what we state in Algorithm 1. That an SRFT with a smaller value of  $r_1$  can be used instead is a consequence of the fact that (Drineas et al., 2012, Lemma 3) is also satisfied by an SRFT matrix with the given  $r_1$ ; this is established in (Tropp, 2011; Boutsidis and Gittens, 2013).

The running time of this algorithm, given in the caption of the algorithm, is roughly  $O(nd \ln d)$  when  $d = \Omega(\ln n)$ . Thus Algorithm 1 generates relative-error approximations to the leverage scores of a tall and skinny matrix  $\mathbf{A}$  in time  $o(nd^2)$ , rather than the  $\Omega(nd^2)$  time that would be required to compute a QR decomposition or a thin SVD of the  $n \times d$  matrix  $\mathbf{A}$ . The basic idea behind Algorithm 1 is as follows. If we had a QR decomposition of  $\mathbf{A}$ , then we could postmultiply  $\mathbf{A}$  by the inverse of the “ $R$ ” matrix to obtain an orthogonal matrix spanning the column space of  $\mathbf{A}$ ; and from this  $n \times d$  orthogonal matrix, we could read off

**Input:**  $\mathbf{A} \in \mathbb{R}^{n \times d}$  (with SVD  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ ), error parameter  $\epsilon \in (0, 1/2]$ .

**Output:**  $\tilde{\ell}_i, i = 1, \dots, n$ , approximations to the leverage scores of  $\mathbf{A}$ .

1. Let  $\mathbf{\Pi}_1 \in \mathbb{R}^{r_1 \times n}$  be an SRFT with

$$r_1 = \Omega(\epsilon^{-2}(\sqrt{d} + \sqrt{\ln n})^2 \ln d)$$

2. Compute  $\mathbf{\Pi}_1\mathbf{A} \in \mathbb{R}^{r_1 \times d}$  and its QR factorization  $\mathbf{\Pi}_1\mathbf{A} = \mathbf{Q}\mathbf{R}$ .

3. Let  $\mathbf{\Pi}_2 \in \mathbb{R}^{d \times r_2}$  be a matrix of i.i.d. standard Gaussian random variables, where

$$r_2 = \Omega(\epsilon^{-2} \ln n).$$

4. Construct the product  $\mathbf{\Omega} = \mathbf{A}\mathbf{R}^{-1}\mathbf{\Pi}_2$ .

5. For  $i = 1, \dots, n$  compute  $\tilde{\ell}_i = \|\Omega_{(i)}\|_2^2$ .

**Algorithm 1:** Algorithm (Drineas et al., 2012, Algorithm 1) for approximating the leverage scores  $\ell_i$  of an  $n \times d$  matrix  $\mathbf{A}$ , where  $n \gg d$ , to within a multiplicative factor of  $1 \pm \epsilon$ . The running time of the algorithm is  $O(nd \ln(\sqrt{d} + \sqrt{\ln n}) + nd\epsilon^{-2} \ln n + d^2\epsilon^{-2}(\sqrt{d} + \sqrt{\ln n})^2 \ln d)$ .

the leverage scores from the Euclidean norms of the rows. Of course, computing the QR decomposition would require  $O(nd^2)$  time. To get around this, Algorithm 1 premultiplies  $\mathbf{A}$  by a structured random projection  $\mathbf{\Pi}_1$ , computes a QR decomposition of  $\mathbf{\Pi}_1\mathbf{A}$ , and postmultiplies  $\mathbf{A}$  by  $\mathbf{R}^{-1}$ , *i.e.*, the inverse of the “ $R$ ” matrix from the QR decomposition of  $\mathbf{\Pi}_1\mathbf{A}$ . Since  $\mathbf{\Pi}_1$  is an SRFT, premultiplying by it takes roughly  $O(nd \ln d)$  time. In addition, note that  $\mathbf{\Pi}_1\mathbf{A}$  needs to be post multiplied by a second random projection in order to compute all of the leverage scores in the allotted time; see (Drineas et al., 2012) for details. This algorithm is simpler than the algorithm in which we are primarily interested that is applicable to square SPSD matrices, but we start with it since it illustrates the basic ideas of how our main algorithm works and since our main algorithm calls it as a subroutine. We note, however, that this algorithm is directly useful for approximating the leverage scores of Linear Kernel matrices  $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ , when  $\mathbf{X}$  is a tall and skinny matrix.

Consider, next, Algorithm 2 (which originally appeared as Algorithm 4 in (Drineas et al., 2012)), which takes as input an *arbitrary*  $n \times d$  matrix  $\mathbf{A}$  and a rank parameter  $k$ , and returns as output a  $1 \pm \epsilon$  approximation to *all* of the statistical leverage scores (relative to the best rank- $k$  approximation) of the input. An important technical point is that the problem of computing the leverage scores of a matrix relative to a low-dimensional space is ill-posed, essentially because the spectral gap between the  $k^{\text{th}}$  and the  $(k+1)^{\text{st}}$  eigenvalues can be small, and thus Algorithm 2 actually computes approximations to the leverage scores of a matrix that is near to  $\mathbf{A}$  in the spectral norm (or the Frobenius norm if  $q = 0$ ). See (Drineas et al., 2012) for details. Basically, this algorithm uses Gaussian sampling to find

**Input:**  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , a rank parameter  $k$ , and an error parameter  $\epsilon \in (0, 1/2]$ .

**Output:**  $\hat{\ell}_i, i = 1, \dots, n$ , approximations to the leverage scores of  $\mathbf{A}$  filtered through its dominant dimension- $k$  subspace.

1. Construct  $\mathbf{\Pi} \in \mathbb{R}^{d \times 2k}$  with i.i.d. standard Gaussian entries.
2. Compute  $\mathbf{B} = (\mathbf{A}\mathbf{A}^T)^q \mathbf{A}\mathbf{\Pi} \in \mathbb{R}^{n \times 2k}$  with

$$q \geq \left\lceil \frac{\ln \left( 1 + \sqrt{\frac{k}{k-1}} + e\sqrt{\frac{2}{k}} \sqrt{\min\{n, d\} - k} \right)}{2 \ln(1 + \epsilon/10) - 1/2} \right\rceil.$$

3. Approximate the leverage scores of  $\mathbf{B}$  by calling Algorithm 1 with inputs  $\mathbf{B}$  and  $\epsilon$ ; let  $\hat{\ell}_i$  for  $i = 1, \dots, n$  be the outputs of Algorithm 1.

**Algorithm 2:** Algorithm (Drineas et al., 2012, Algorithm 4) for approximating the leverage scores (relative to the best rank- $k$  approximation to  $\mathbf{A}$ ) of a general  $n \times d$  matrix  $\mathbf{A}$  with those of a matrix that is close by in the spectral norm (or the Frobenius norm if  $q = 0$ ). This algorithm runs in time  $O(ndkq) + T_1$ , where  $T_1$  is the running time of Algorithm 1.

a matrix close to  $\mathbf{A}$  in the Frobenius norm or spectral norm, and then it approximates the leverage scores of this matrix by using Algorithm 1 on the smaller, very rectangular matrix  $\mathbf{B}$ . When  $\mathbf{A}$  is square, as in our applications, Algorithm 2 is typically more costly than direct computation of the leverage scores, at least for dense matrices (but it does have the advantage that the number of iterations is bounded, independent of properties of the matrix, which is not true for typical iterative methods to compute low-rank approximations).

Of greater practical interest is Algorithm 3, which is a modification of Algorithm 2 in which the Gaussian random projection is replaced with an SRFT. That is, Algorithm 3 uses an SRFT projection to find a matrix close by to  $\mathbf{A}$  in the Frobenius norm or spectral norm (depending on the value of  $q$ ), and then it exactly computes the leverage scores of this matrix. This improves the running time to  $O(n^2 \ln(\sqrt{k} + \sqrt{\ln n}) + n^2(\sqrt{k} + \sqrt{\ln n})^2 \ln(k)q + n(\sqrt{k} + \sqrt{\ln n})^4 \ln^2(k))$ , which is  $o(n^2k)$  when  $q = 0$ . Thus an important point for Algorithm 3 (as well as for Algorithm 2) is the parameter  $q$  which describes the number of iterations. For  $q = 0$  iterations, we get an inexpensive Frobenius norm approximation; while for higher  $q$ , we get better spectral norm approximations that are more expensive.<sup>3</sup> This flexibility is of interest, as one may want to approximate the actual leverage scores accurately or one may simply want to find crude approximations useful for obtaining SPSP sketches with low reconstruction error.

3. Observe that since  $\mathbf{A}$  is rectangular in Algorithms 2 and 3, we approximate the leverage scores of  $\mathbf{A}$  with those of  $\mathbf{B} = (\mathbf{A}\mathbf{A}^T)^q \mathbf{A}\mathbf{\Pi}$ ; in particular the case  $q = 0$  corresponds to taking  $\mathbf{B} = \mathbf{A}\mathbf{\Pi}$ . By way of contrast, when we use the power method to construct sketches of an SPSP matrix, we take  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$ , so the case  $q = 1$  corresponds to  $\mathbf{C} = \mathbf{A}\mathbf{S}$ .

**Input:**  $\mathbf{A} \in \mathbb{R}^{n \times d}$ , a rank parameter  $k$ , and an iteration parameter  $q$ .

**Output:**  $\hat{\ell}_i, i \in \{1, \dots, n\}$ , approximations to the leverage scores of  $\mathbf{A}$  filtered through its dominant dimension- $k$  subspace.

1. Construct an SRHT matrix  $\mathbf{\Pi} \in \mathbb{R}^{d \times r}$ , where

$$r \geq \left\lceil 36\epsilon^{-2}[\sqrt{k} + \sqrt{8\ln(kd)}]^2 \ln(k) \right\rceil.$$

2. Compute  $\mathbf{B} = (\mathbf{A}\mathbf{A}^T)^q \mathbf{A}\mathbf{\Pi} \in \mathbb{R}^{n \times r}$ , where  $q \geq 0$  is an integer.
3. Return the exact leverage scores of  $\mathbf{B}$ .

**Algorithm 3:** Algorithm for approximating the leverage scores (relative to the best rank- $k$  approximation to  $\mathbf{A}$ ) of a general  $n \times d$  matrix  $\mathbf{A}$  with those of a matrix that is close by in the spectral norm. This is a modified version of Algorithm 2, in which the random projection is implemented with an SRFT rather than a Gaussian random matrix, and where the number of “iterations”  $q$  is prespecified. This algorithm runs in time  $O(nd \ln r + ndr q + nr^2)$  since  $\mathbf{A}\mathbf{\Pi}$  can be computed in time  $O(nd \ln r)$ .

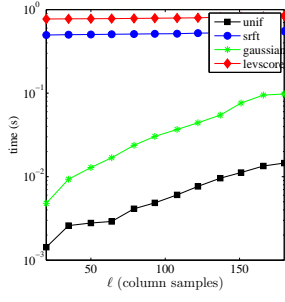
Finally, note that although choosing the number of iterations  $q$  as we did in Algorithm 2 is convenient for worst-case analysis, as a practical implementational matter it is easier either to choose  $q$  based on spectral gap information revealed during the running of the algorithm or to prespecify  $q$  to be a small integer, *e.g.*, 2 or 3, before the algorithm runs. Both of these have an interpretation of accelerating the rate of decay of the spectrum with a power iteration, but they behave somewhat differently due to the different stopping conditions. Below, we consider both variants.

### 3.5.2 RUNNING TIME COMPARISONS

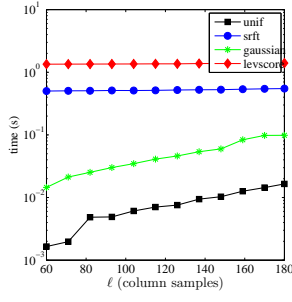
Here, we describe the performances of the various random sampling and random projection low-rank sketches considered in Section 3.4 in terms of their running time, where the method that involves using the leverage scores to construct the importance sampling distribution is implemented both by computing the leverage scores “exactly” by calling a truncated SVD, as a black box, as well as computing them approximately by using one of several versions of Algorithm 3. Our running time results are presented in Figure 6 and Figure 7.

We start with the results described in Figure 6, which shows the running times, as a function of  $\ell$ , for the low-rank approximations described in Section 3.4: *i.e.*, for column sampling uniformly at random without replacement; for column sampling according to the exact nonuniform leverage score probabilities; and for sketching using Gaussian and SRFT mixtures of the columns. Several observations are worth making about the results presented in this figure.

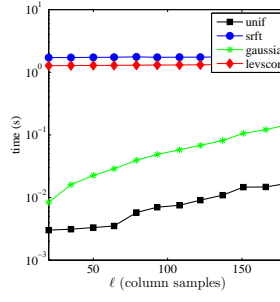
REVISITING THE NYSTRÖM METHOD



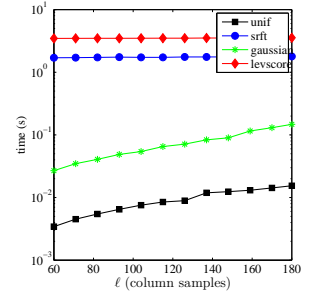
(a) GR,  $k = 20$



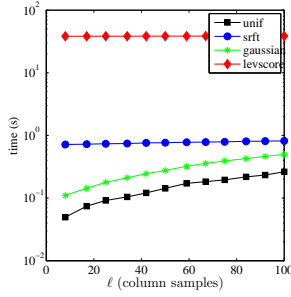
(b) GR,  $k = 60$



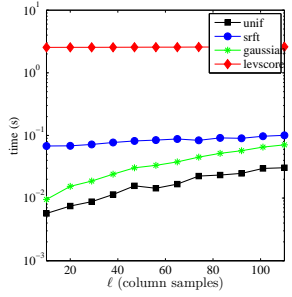
(c) HEP,  $k = 20$



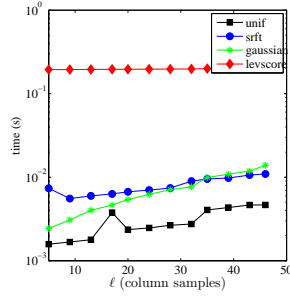
(d) HEP,  $k = 60$



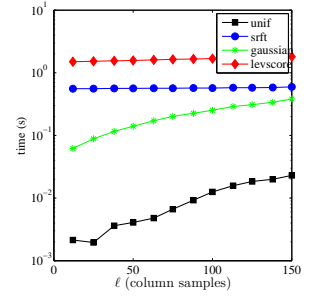
(e) Dexter,  $k = 8$



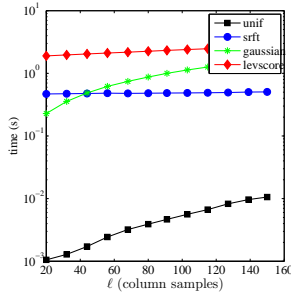
(f) Protein,  $k = 10$



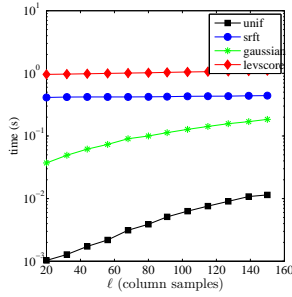
(g) SNPs,  $k = 5$



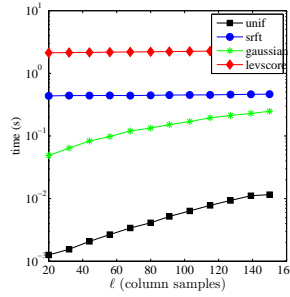
(h) Gisette,  $k = 12$



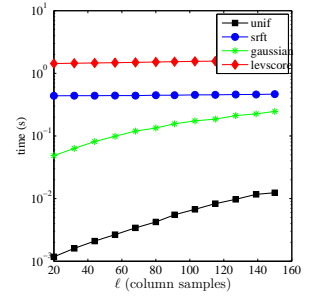
(i) AbaloneD,  $\sigma = .15, k = 20$



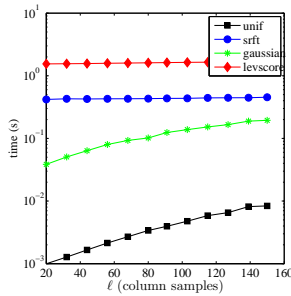
(j) AbaloneD,  $\sigma = 1, k = 20$



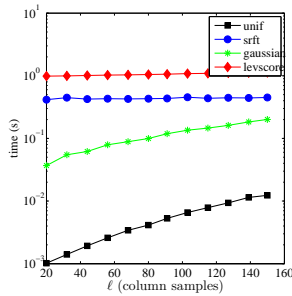
(k) WineD,  $\sigma = 1, k = 20$



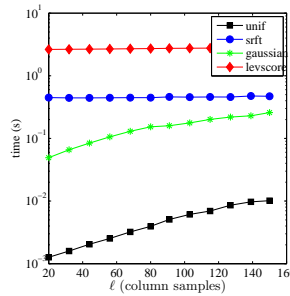
(l) WineD,  $\sigma = 2.1, k = 20$



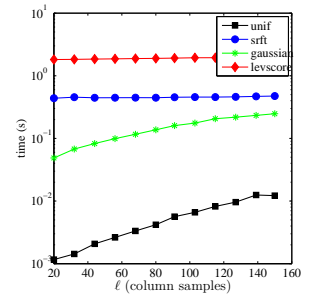
(m) AbaloneS,  $\sigma = .15, k = 20$



(n) AbaloneS,  $\sigma = 1, k = 20$



(o) WineS,  $\sigma = 1, k = 20$



(p) WineS,  $\sigma = 2.1, k = 20$

Figure 6: The times required to compute SPSD sketches, as a function of the number of column samples  $\ell$  for several data sets and two choices of the rank parameter  $k$ .

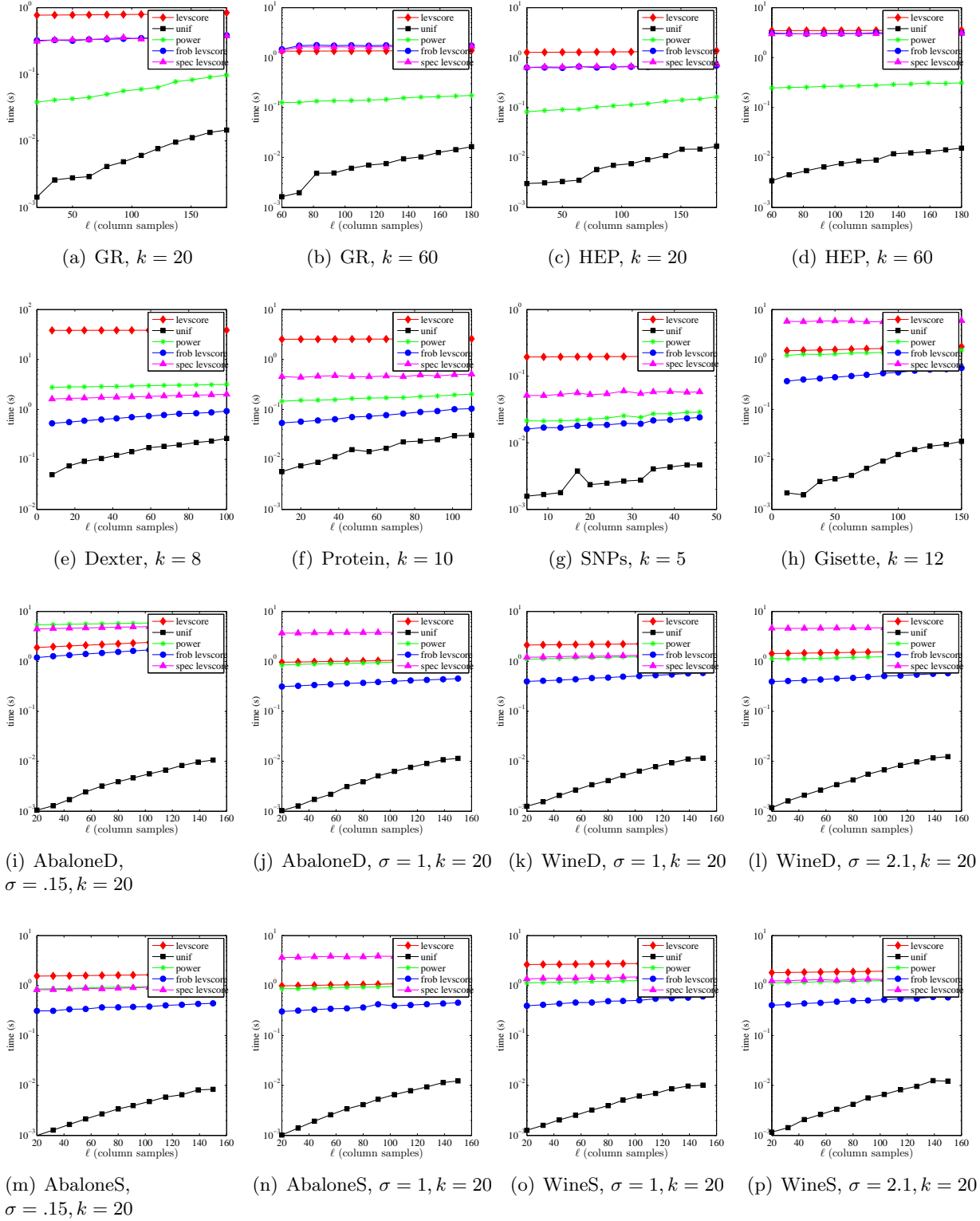


Figure 7: The times required to compute approximate leverage score-based SPSD sketches, as a function of the number of column samples  $\ell$  for several data sets.

- Uniform sampling is always less expensive and typically much less expensive than the other methods, while (with one minor exception) sampling according to the *exact* leverage scores is always the most expensive method.
- For most matrices, using the SRFT is nearly as expensive as exact leverage score sampling. This is most true for the very sparse graph Laplacian Kernels, largely since the SRFT does not respect sparsity. The main exception to this is for the dense and relatively well-behaved Linear Kernels, where especially for large values of  $\ell$  the SRFT is quite fast and usually not too much more expensive than uniform sampling.
- The “fast Fourier” methods underlying the SRFT can take advantage of the structure of the Linear Kernels to yield algorithms that are similar to Gaussian projections and much better than exact leverage score computation. Note that the reason that SRFT is worse than Gaussians here is that the matrices we are considering are *not* extremely large, and we are not considering very large values of the rank parameter. Extending in both those directions leads to Gaussian projections being slower than SRFT, as the trends in the figures clearly indicate.
- Gaussian projections are not too much slower than uniform sampling for the extremely sparse Laplacian Kernels—this is due to the sparsity of the Laplacian Kernels, since Gaussian projections can take advantage of fast matrix-vector multiplies, while the SRFT-based scheme cannot—but this advantage is lost for the (denser) Sparse RBF Kernels, to the extent that there is little running time improvement relative to the Dense RBF Kernels. In addition, Gaussian projections are relatively slower, when compared to the SRFT and uniform sampling, for the Dense RBF Kernels than for the Linear Kernels, although both of those data sets are maximally dense.

We next turn to the results described in Figure 7, which shows the running times, as a function of  $\ell$ , for several variants of approximate leverage-based sampling. For ease of comparison, the timings for uniform sampling (“unif”) and exact leverage score sampling (“levscore”) are depicted in Figure 7 using the same shading as used in Figure 6. In addition to these two baselines, Figure 7 shows running time results for the following three variants of approximate leverage score sampling: “frob levscore” (which is Algorithm 3 with  $q = 0$  and  $r = 2k$ ); “spec levscore” (Algorithm 3 with  $q = 4$  and  $r = 2k$ ); and “power”. The “power” scheme is a version of Algorithm 3 where  $r = k$  and  $q$  is determined by monitoring the convergence of the leverage scores of  $\mathbf{A}^{2q+1}\mathbf{\Pi}$  and terminating when the change in the leverage scores between iterations, as measured in the infinity norm, is smaller than  $10^{-2}$ . This is simply a version of subspace iteration with a convergence criterion appropriate for the task at hand. Since “frob levscore” requires one application of an SRFT, its timing results are depicted using the same shade as the SRFT timing results in Figure 6. (There are no other correspondences between the shadings in the two figures.) Several observations are worth making about the results presented in this figure.

- These approximate leverage score-based algorithms can be orders of magnitude faster than exact leverage score computation; but, especially for “spec levscore” when  $q$  is not prespecified to be 2 or 3, they can even be somewhat slower. Exactly which is

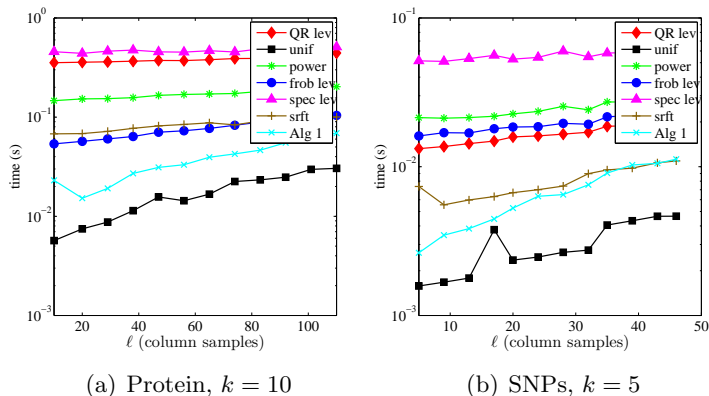


Figure 8: The running time of SPSD sketches computed using Algorithm 1 compared with that of other approximate leverage score-based SPSD sketches, as a function of the number of column samples  $\ell$  for two Linear Kernel datasets. The parameters in Algorithm 1 were taken to be  $r_1 = \epsilon^{-2} \ln(d\delta^{-1})(\sqrt{d} + \sqrt{\ln(n\delta^{-1})})^2$  and  $r_2 = \epsilon^{-2}(\ln n + \ln \delta^{-1})$  with  $\epsilon = 1$  and  $\delta = 1/10$ .

the case depends upon the properties of the matrix and the parameters used in the approximation algorithm, including especially the number of power iterations.

- The “frob levscore” approximation method has running time comparable to the running time of the SRFT, which is expected, given that the computation of the SRFT is the theoretical bottleneck for the running time of the “frob levscore” algorithm. In particular, for larger values of  $\ell$  for Linear Kernels, “frob levscore” is not much slower than uniform sampling.
- The “spec levscore” and “power” approximations with  $q > 0$  are more expensive than the  $q = 0$  “frob lev” approximation, which is a result of the relatively-expensive matrix-matrix multiplication. For the Linear Kernels, both are much better than the exact leverage score computation, and for most other data at least “power” is somewhat less expensive than the exact leverage score computation. For example, this is particularly true for the Laplacian Kernels.

Recall that the cost associated with these SPSD sketches is two-fold: first, the cost to construct the sample—by sampling columns uniformly at random, by computing a nonuniform importance sampling distribution, or by performing a random projection to uniformize the leverage scores; and second, the cost to construct the low-rank approximation from the sample. For uniform sampling, the latter step dominates the cost, while for more sophisticated methods the former step typically dominates the cost. The approximate leverage score sampling methods are still sufficiently expensive that the cost of computing the sampling probabilities still dominates the cost to construct the low-rank approximation.

Finally, Algorithm 1 can be used to approximate quickly the leverage scores of matrices of the form  $\mathbf{A} = \mathbf{X}\mathbf{X}^T$ , when  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is a rectangular matrix of sufficient aspect ratio, and



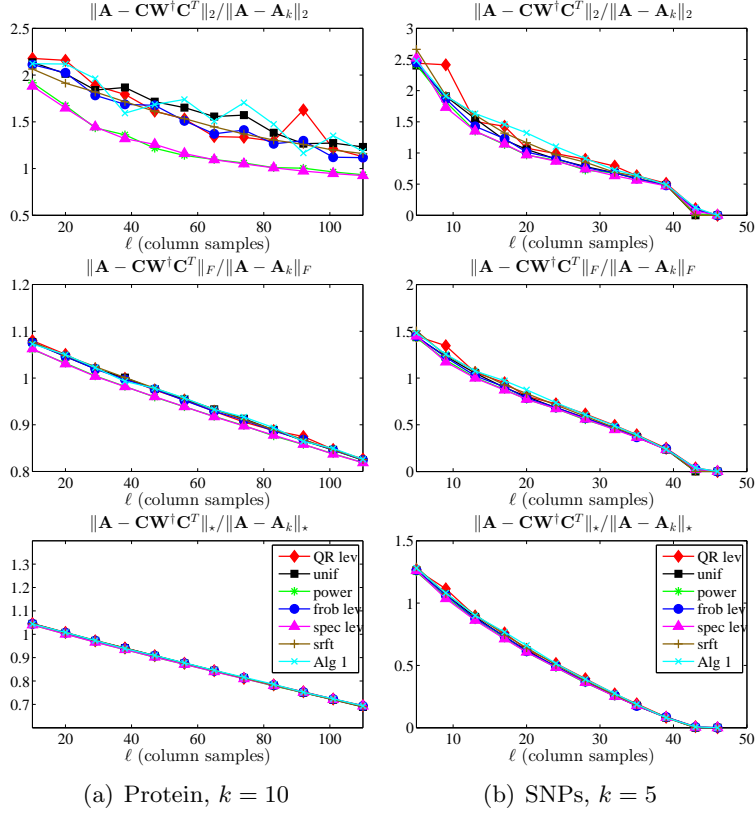


Figure 9: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of SPSD sketches computed using Algorithm 1 compared with those of other approximate leverage score-based sketching schemes, as a function of the number of column samples  $\ell$ , for two Linear Kernel data sets. The parameters in Algorithm 1 were taken to be  $r_1 = \epsilon^{-2} \ln(d\delta^{-1})(\sqrt{d} + \sqrt{\ln(n\delta^{-1})})^2$  and  $r_2 = \epsilon^{-2}(\ln n + \ln \delta^{-1})$  with  $\epsilon = 1$  and  $\delta = 1/10$ .

in such cases it is faster than Algorithm 3. Specifically, for the first dimensional reduction step in Algorithm 1 to be beneficial (*i.e.*, to ensure  $r_1 < n$ ), the condition  $n = \Omega(d \ln d)$  is necessary; for the second dimensional reduction step to be beneficial (*i.e.*, to ensure  $r_2 < d$ ), the condition  $d = \Omega(\ln n)$  must be satisfied. Figure 8 summarizes our main results for the run time of Algorithm 1 applied to rectangular matrices with  $n \gg d$ . Among other things, Figure 8 illustrates, using the Linear Kernel datasets Protein and SNPs (which satisfy these constraints), two points.

- Most importantly, the running time of Algorithm 1 on these rectangular matrices is faster than performing a QR decomposition on  $\mathbf{A}$  and is comparable to applying a SRFT to  $\mathbf{A}$ . This is expected, since the running time bottleneck for Algorithm 1 is the application of the SRFT.
- In addition, the running time of Algorithm 1 is significantly faster than the other approximate leverage score algorithms. This too is expected, since these other algorithms are applied to  $\mathbf{A}$  and ignore the rectangular structure of  $\mathbf{X}$ .

Figure 9 shows that these improved running time gains for Algorithm 1 can come at the cost of a slight loss in the reconstruction accuracy (relative to the exact computation of the leverage scores) of the low-rank approximations; the accuracy of the other approximate leverage score algorithms is discussed in the following subsection.

### 3.5.3 RECONSTRUCTION ACCURACY RESULTS

Here, we describe the performances of the various low-rank approximations that use approximate leverage scores in terms of reconstruction accuracy for the data sets described in Section 3.1. The results are presented in Figure 10 through Figure 14. The setup for these results parallels that for the low-rank approximation results described in Section 3.4, and these figures parallel Figure 1 through Figure 5. To provide a baseline for the comparison, we also plot the previous reconstruction errors for sampling with the exact leverage scores as well as the uniform column sampling sketch. Several observations are worth making about the results presented in these figures.

For Laplacian Kernels, “frob levscore” is only slightly better than uniform sampling, while “power” and “spec levscore” are substantially better than uniform sampling; all of those methods also lead to even *better* reconstruction results than using the exact leverage scores (suggesting that some form of implicit regularization is taking place): the reconstruction quality is higher for a given  $\ell$  and, also, using approximate leverage scores does not lead to the saturation effect observed when using the exact leverage scores. For the Linear Kernels, all the methods perform similarly. For both the dense and the sparse RBF data sets, the approximate leverage score algorithms tend to parallel the exact leverage score algorithm, and they are not substantially better. In particular, both “power” and “spec levscore” tend to saturate when the exact method saturates, but in those cases “frob levscore” tends not to saturate.

Note that the difference between different approximate leverage score algorithms often corresponds to a difference in the spectral gaps of the corresponding matrices. From Table 5, if we fix  $k$  and use the approximate leverage scores filtered through rank  $k$  to form a Nyström approximation to  $\mathbf{A}$ , the accuracy of that approximation has a strong dependence on the

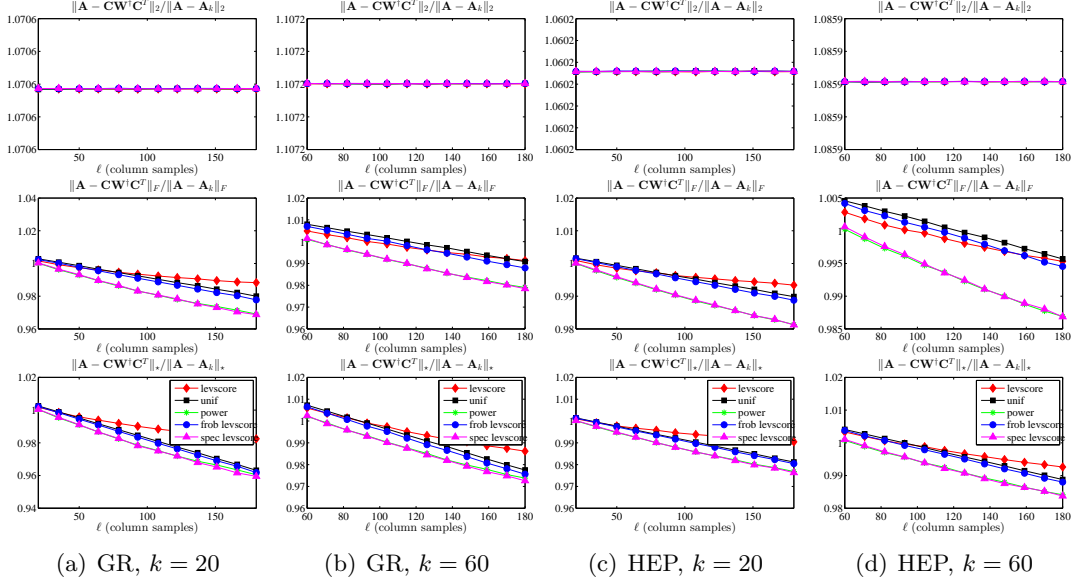


Figure 10: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several approximate leverage score-based PSD sketches, as a function of the number of column samples  $\ell$ , for the GR and HEP Laplacian data sets, with two choices of the rank parameter  $k$ .

spectral gap of  $\mathbf{A}$  at rank  $k$ , as measured by  $\frac{\lambda_k}{\lambda_{k+1}}$ . In general, the larger the spectral gap, the more accurate the approximation. This phenomena can also be understood in terms of the convergence of the approximate leverage scores: the approximation algorithms (Algorithm 2 and Algorithm 3) are essentially truncated versions of the subspace iteration method for computing the top  $k$  eigenvectors of  $\mathbf{A}$ . It is a classical result that the spectral gap determines the rate of convergence of the subspace iteration process to the desired eigenvectors: the larger it is, the fewer iterations of the process are required to get accurate approximations of the top eigenvectors. It follows immediately that the larger the spectral gap, the more accurate the approximate leverage scores generated by these approximation algorithms are. Our empirical results illustrate the complexities and subtle consequences of these properties in realistic machine learning applications of even modestly-large size.

### 3.5.4 SUMMARY OF LEVERAGE SCORE APPROXIMATION ALGORITHMS

Before proceeding, there are several summary observations that we can make about the running time and reconstruction quality of approximate leverage score sampling algorithms for the data sets we have considered.

- The running time of computing the exact leverage scores is generally much worse than that of uniform sampling and both SRFT-based and Gaussian-based random projection methods.

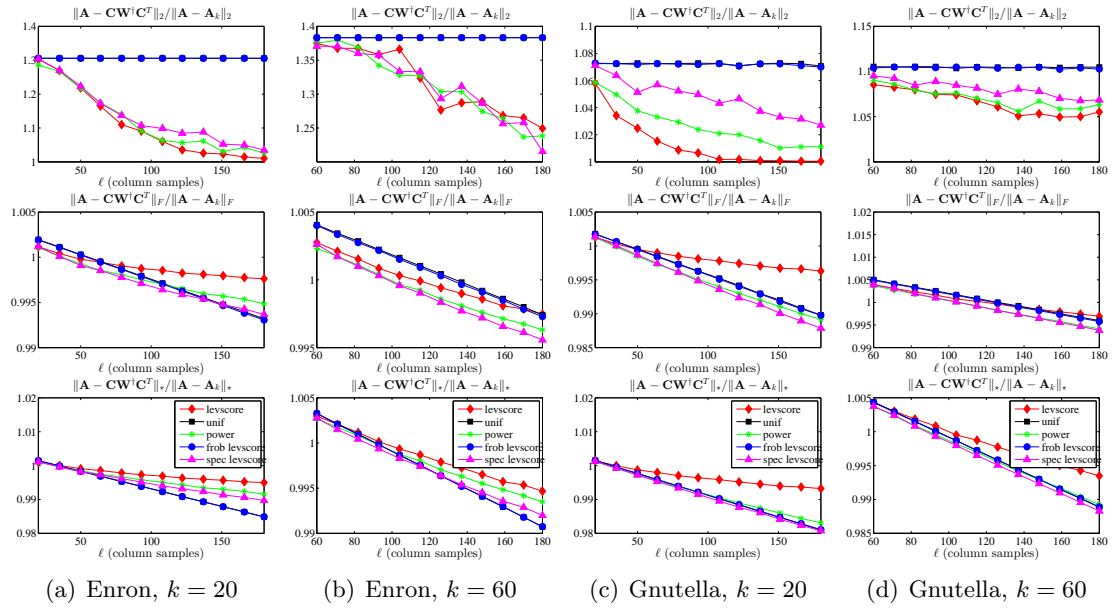


Figure 11: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several approximate leverage score-based PSD sketches, as a function of the number of column samples  $\ell$ , for the Enron and Gnutella Laplacian data sets, with two choices of the rank parameter  $k$ .

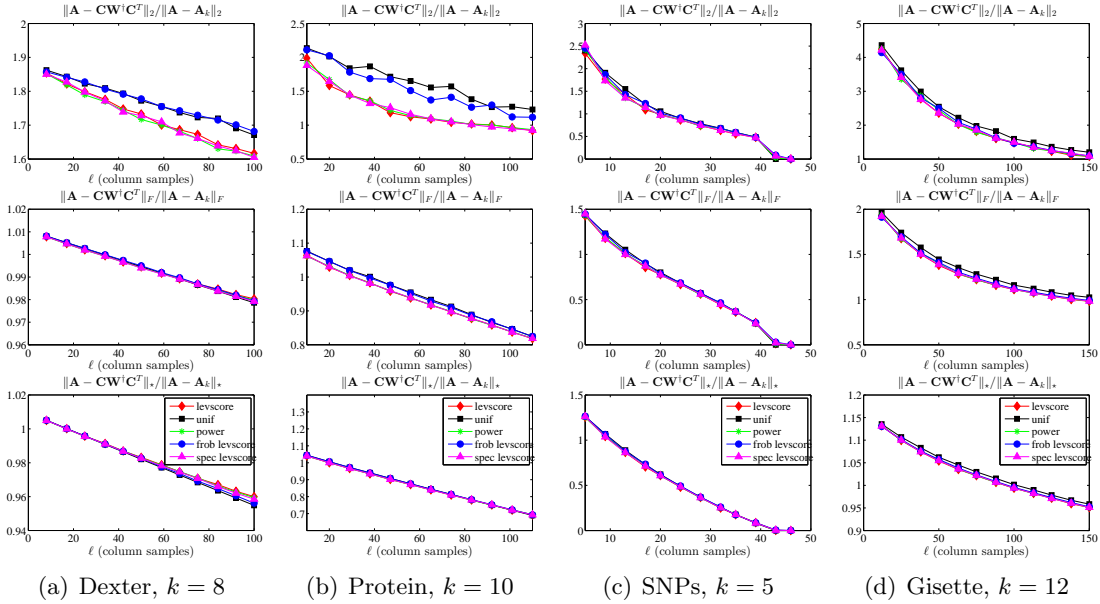


Figure 12: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several approximate leverage score-based SPSD sketches, as a function of the number of column samples  $\ell$ , for the Linear Kernel data sets.

- The running time of computing approximations to the leverage scores can, with appropriate choice of parameters, be much faster than the exact computation of the leverage scores; and, especially for “frob levscore,” can be comparable to the running time of the random projection (SRFT or Gaussian) used in the leverage score approximation algorithm. For the methods that involve  $q > 0$  iterations to compute stronger approximations to the leverage scores, the running time can vary considerably depending on details of the stopping condition.
- The leverage scores computed by the “frob levscore” procedure are typically very different than the “exact” leverage scores, but they are leverage scores for a low-rank space that is near the best rank- $k$  approximation to the matrix. This is often sufficient for good low-rank approximation.
- The approximate leverage scores computed from “power” and “spec levscore” approach those of the exact leverage scores, as  $q$  is increased; and they obtain reconstruction accuracy that is no worse, and in many cases is better, than that obtained by the exact leverage scores. This suggests that, by not fitting exactly to the empirical statistical leverage scores, we are observing a form of implicit regularization.
- The running time of Algorithm 1, when applied to “tall” matrices for which  $n \gg d$ , is faster than the running time of performing a QR decomposition of the matrix  $\mathbf{A}$ ; and it is comparable to the running time of applying a random projection to  $\mathbf{A}$  (which is the computational bottleneck of applying Algorithm 1). Thus, in particular, one could

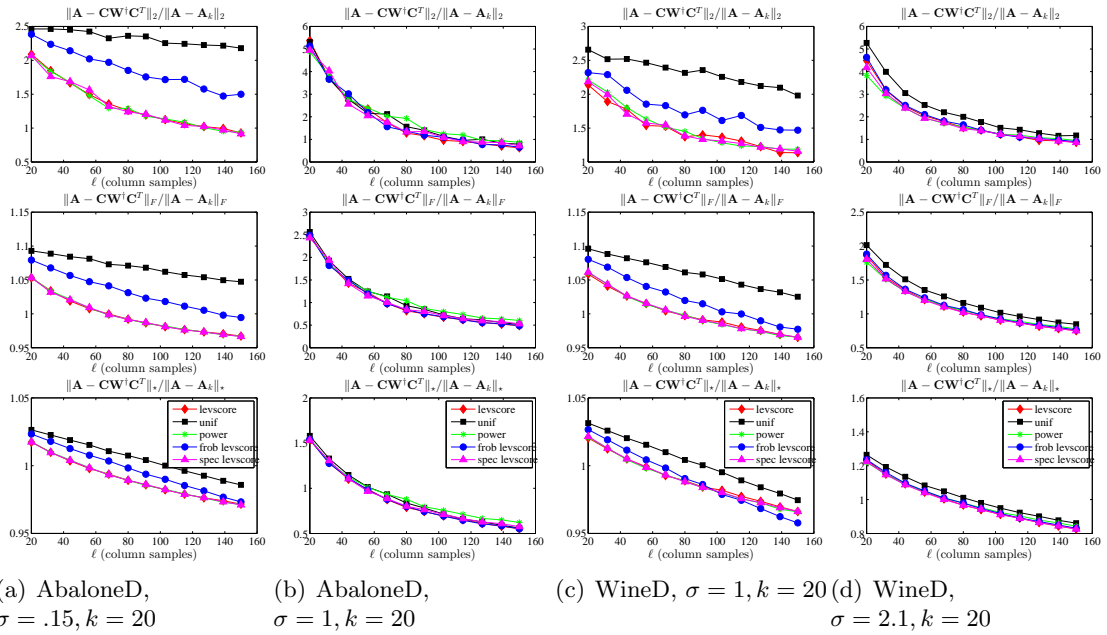


Figure 13: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several approximate leverage score-based PSD sketches, as a function of the number of column samples  $\ell$ , for several dense RBF data sets.

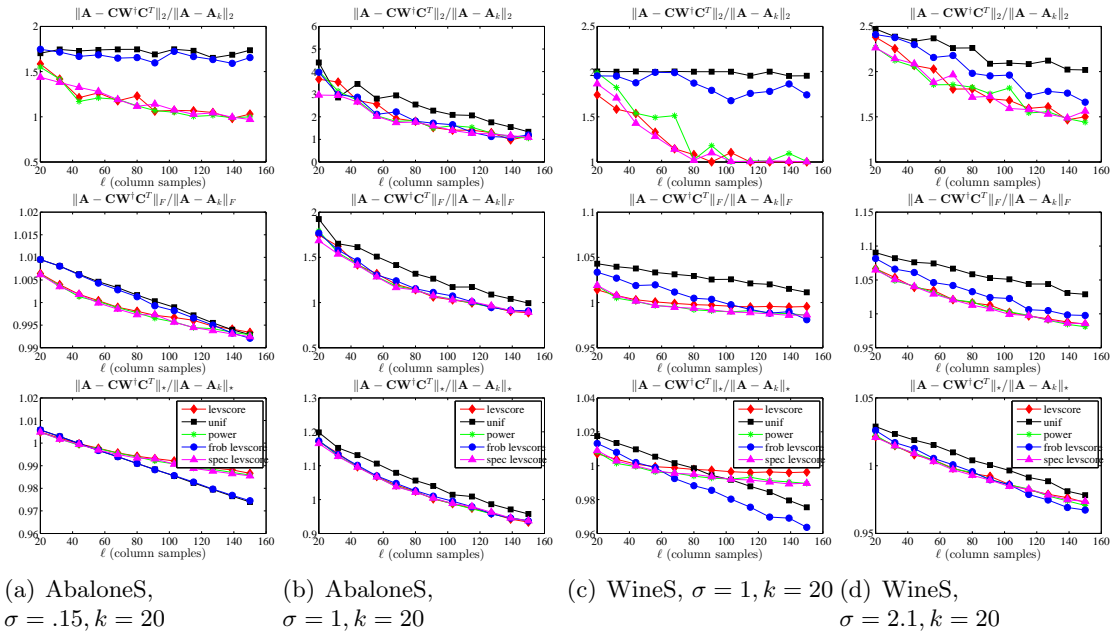


Figure 14: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several approximate leverage score-based PSD sketches, as a function of the number of column samples  $\ell$ , for several sparse RBF data sets.

use this algorithm to compute approximations to the leverage scores to obtain a sketch that provides a relative-error approximation to a least-squares problem involving  $\mathbf{A}$  (Drineas et al., 2008, 2010; Mahoney, 2011); or one could use the sketch thereby obtained as a preconditioner to an iterative method to solve the least-squares problem, in a manner analogous to how Blendenpik or LSRN do so with a random projection (Avron et al., 2010; Meng et al., 2014).

Previous work has showed that one can implement random projection algorithms to provide low-rank approximations with error comparable to that of the SVD in less time than state-of-the-art Krylov solvers and other “exact” numerical methods (Halko et al., 2011; Mahoney, 2011). Our empirical results show that these random projection algorithms can be used in two complementary ways to approximate SPSP matrices of interest in machine learning: first, they can be used directly to compute a projection-based low-rank approximation; and second, they can be used to compute approximations to the leverage scores, which can be used to compute a sampling-based low-rank approximation. With the right choice of parameters, the two complementary approaches have roughly comparable running times, and neither one dominates the other in terms of reconstruction accuracy.

### 3.6 Projection-based Sketches

Finally, for completeness, we consider the performance of the two projection-based SPSP sketches proposed by Halko et al. (2011), and we show how they perform when compared with the sketches we have considered. Recall that the idea of these sketches is to construct low-rank approximations by forming an approximate basis  $\mathbf{Q}$  for the top eigenspace of  $\mathbf{A}$  and then restricting  $\mathbf{A}$  to that eigenspace. In more detail, given a sketching matrix  $\mathbf{S}$ , form the matrix  $\mathbf{Y} = \mathbf{A}\mathbf{S}$  and take the QR decomposition of  $\mathbf{Y}$  to obtain  $\mathbf{Q}$ , a matrix with orthonormal columns. The first sketch, which we eponymously refer to as the *pinched* sketch, is simply  $\mathbf{A}$  pinched to the space spanned by  $\mathbf{Q}$ :

$$\mathbf{Q}(\mathbf{Q}^T \mathbf{A} \mathbf{Q}) \mathbf{Q}^T. \quad (8)$$

The second sketch, which we refer to as the *prolonged* sketch, is

$$\mathbf{A} \mathbf{Q} (\mathbf{Q}^T \mathbf{A} \mathbf{Q})^\dagger \mathbf{Q}^T \mathbf{A}. \quad (9)$$

It is clear that the prolonged sketch can be constructed using our SPSP Sketching Model by taking  $\mathbf{Q}$  as the sketching matrix. In fact, a stronger statement can be made. As stated in Lemma 1 below, it is the case, for any sketching matrix  $\mathbf{X}$ , that when  $\mathbf{C} = \mathbf{A}\mathbf{X}$  and  $\mathbf{W} = \mathbf{X}^T \mathbf{A}\mathbf{X}$ ,

$$\mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T = \mathbf{A}^{1/2} \mathbf{P}_{\mathbf{A}^{1/2} \mathbf{X}} \mathbf{A}^{1/2}.$$

By considering the sketching matrix  $\mathbf{X} = \mathbf{A}^1 \mathbf{S}$ , we see that in fact the prolonged sketch is exactly the sketch obtained by applying the power method with  $q = 1$ :

$$\begin{aligned} \mathbf{A} \mathbf{Q} (\mathbf{Q}^T \mathbf{A} \mathbf{Q})^\dagger \mathbf{Q}^T \mathbf{A} &= \mathbf{A}^{1/2} \mathbf{P}_{\mathbf{A}^{1/2} \mathbf{Q}} \mathbf{A}^{1/2} \\ &= \mathbf{A}^{1/2} \mathbf{P}_{\mathbf{A}^{1/2} (\mathbf{A}\mathbf{S})} \mathbf{A}^{1/2} \\ &= \mathbf{A}^2 \mathbf{S} (\mathbf{S}^T \mathbf{A}^3 \mathbf{S})^\dagger \mathbf{S}^T \mathbf{A}^2 \\ &= \mathbf{A} \mathbf{X} (\mathbf{X}^T \mathbf{A} \mathbf{X})^\dagger \mathbf{X}^T \mathbf{A}. \end{aligned}$$



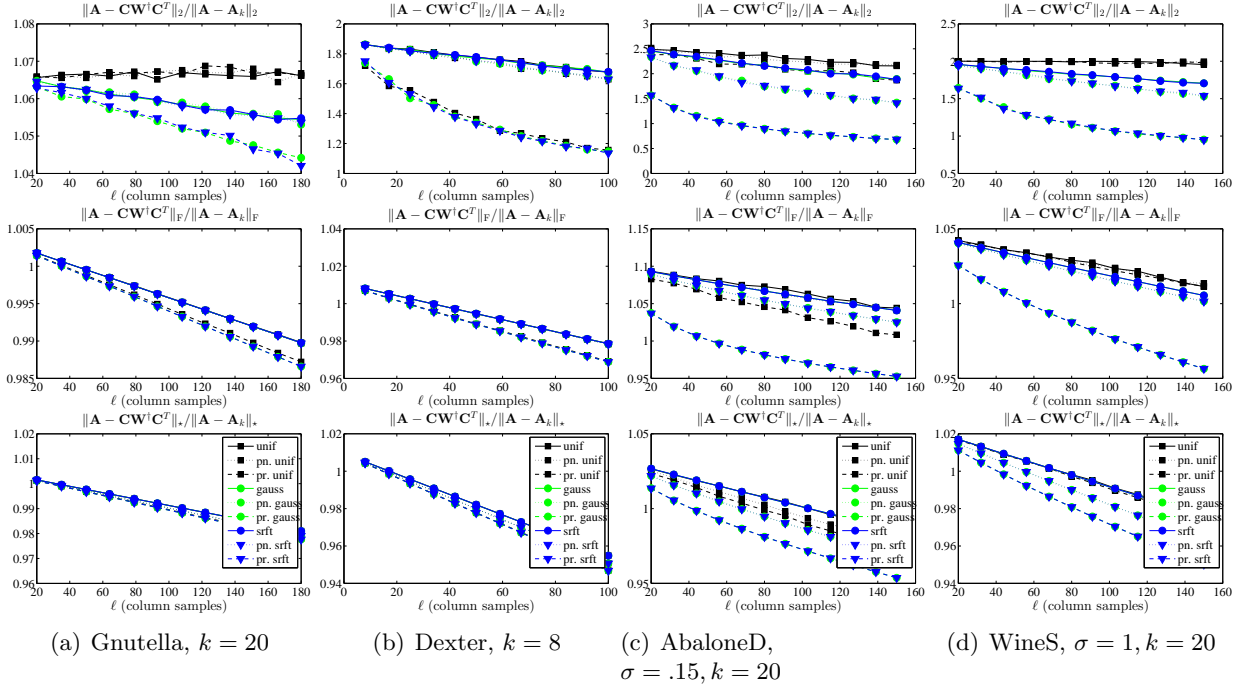


Figure 15: The spectral, Frobenius, and trace norm errors (top to bottom, respectively, in each subfigure) of several SPSD sketches, including the pinched and prolonged sketches, as a function of the number of column samples  $\ell$ , for several datasets. Pinched and prolonged sketches, respectively indicated by “pn.” and “pr.,” are defined in Equations (8) and (9).

It follows that the bounds we provide in Section 4 on the performance of sketches obtained using the power method pertain also to prolonged sketches.

In Figure 15, we compare the empirical performances of several of the SPSD sketches considered earlier with their pinched and prolonged variants. Specifically, we plot the errors of pinched and prolonged sketches for several choices of sketching matrices—corresponding to uniform column sampling, Gaussian column mixtures, and SRFT-based column mixtures—along with the errors of non-pinched, non-prolonged sketches constructed using the same choices of  $\mathbf{S}$ . In the interest of brevity, we provide results only for several of the datasets listed in Table 4.

Some trends are clear from Figure 15.

- In the spectral norm, the prolonged sketches are considerably more accurate than the pinched and standard sketches for all the datasets considered. Without exception, the prolonged Gaussian and SRFT column-mixture sketches are the most accurate in the spectral norm, of all the sketches considered. Only in the case of the Dexter Linear Kernel is the prolonged uniformly column-sampled sketch nearly as accurate in the spectral norm as the prolonged Gaussian and SRFT sketches. To a lesser extent, the prolonged sketches are also more accurate in the Frobenius and trace norms than

the other sketches considered. The increased Frobenius and trace norm accuracy is particularly notable for the two RBF Kernel datasets; again, the prolonged Gaussian and SRFT sketches are considerably more accurate than the prolonged uniformly column-sampled sketches.

- After the prolonged sketches, the pinched Gaussian and SRFT column-mixture sketches exhibit the least spectral, Frobenius, and trace norm errors. Again, however, we see that the pinched uniformly column-sampled sketches are considerably less accurate than the pinched Gaussian and SRFT column-mixture sketches. Particularly in the spectral and Frobenius norms, the pinched uniformly column-sampled sketches are not any more accurate than the basic uniformly column-sampled sketches.

From these considerations, it seems evident that the benefits of pinched and prolonged sketches are most prominent when the spectral norm is the error metric, or when the dataset is an RBF Kernel. In particular, pinched and prolonged sketches are not significantly more accurate (than the sketches considered in the previous subsections) in the Frobenius and trace norms for any of the datasets considered.

It is also evident from Figure 15 that the pinched sketches often have a much slighter increase in accuracy over the basic sketches than do the prolonged sketches. To understand why the pinched sketches are less accurate than the prolonged sketches, observe that the pinched sketches satisfy

$$\begin{aligned} \mathbf{Q}(\mathbf{Q}^T \mathbf{A} \mathbf{Q}) \mathbf{Q}^T &= \mathbf{P}_{\mathbf{A}\mathbf{S}} \mathbf{A} \mathbf{P}_{\mathbf{A}\mathbf{S}} \\ &= (\mathbf{P}_{\mathbf{A}\mathbf{S}} \mathbf{A}^{1/2}) (\mathbf{A}^{1/2} \mathbf{P}_{\mathbf{A}\mathbf{S}}), \end{aligned}$$

while, as noted above, the prolonged sketches can be written in the form

$$\mathbf{A} \mathbf{Q} (\mathbf{Q}^T \mathbf{A} \mathbf{Q})^\dagger \mathbf{Q}^T \mathbf{A} = (\mathbf{A}^{1/2} \mathbf{P}_{\mathbf{A}^{3/2}\mathbf{S}}) (\mathbf{P}_{\mathbf{A}^{3/2}\mathbf{S}} \mathbf{A}^{1/2}).$$

Thus, pinched and prolonged sketches approximate the square root of  $\mathbf{A}$  by projecting, respectively, onto the ranges of  $\mathbf{A}\mathbf{S}$  and  $\mathbf{A}^{3/2}\mathbf{S}$ . The spectral decay present in  $\mathbf{A}$  is increased when  $\mathbf{A}$  is raised to a power larger than one; consequently, the range of  $\mathbf{A}^{3/2}\mathbf{S}$  is more biased towards the top  $k$ -dimensional invariant subspace of  $\mathbf{A}$  than is the range of  $\mathbf{A}\mathbf{S}$ . It follows that the approximate square root used to construct the prolonged sketches more accurately captures the top  $k$ -dimensional subspace of  $\mathbf{A}$  than does that used to construct the pinched sketches.

#### 4. Theoretical Aspects of SPSD Low-rank Approximation

In this section, we present our main theoretical results, which consist of a suite of bounds on the quality of low-rank approximation under several different sketching methods. As mentioned above, these were motivated by our empirical observation that *all* of the sampling and projection methods we considered perform *much* better on the SPSD matrices we considered than previous worst-case bounds (*e.g.*, Drineas and Mahoney, 2005; Kumar et al., 2012; Gittens, 2012) would suggest. We start in Section 4.1 with deterministic structural conditions for the spectral, Frobenius, and trace norms. In Section 4.2, we use these results to provide our bounds for several random sampling and random projection procedures.

#### 4.1 Deterministic Error Bounds for Low-rank SPSD Approximation

In this section, we present three theorems that provide error bounds for the spectral, Frobenius, and trace norm approximation errors under the SPSD Sketching Model of Section 2.2. These bounds hold for *any, e.g.*, deterministic or randomized, sketching matrix  $\mathbf{S}$ . Thus, *e.g.*, one could use them to check, in an *a posteriori* manner, the quality of a sketching method for which one cannot establish an *a priori* bound. Rather than doing this, we use these results (in Section 4.2 below) to derive *a priori* bounds for when the sketching operation consists of common random sampling and random projection algorithms. We note that the bounds can be interpreted geometrically in terms of the angles between the subspace spanned by the sampling matrix  $\mathbf{S}$  and the dominant eigenspaces of  $\mathbf{A}$ ; we refer the interested reader to the technical report (Gittens and Mahoney, 2013) for details.

Our results are based on the fact that approximations which satisfy our SPSD Sketching Model can be written in terms of a projection onto a subspace of the range of the square root of the matrix being approximated. The following fact appears in the proof of (Gittens, 2012, Proposition 1).

**Lemma 1** *Let  $\mathbf{A}$  be an SPSD matrix and  $\mathbf{S}$  be a conformal sketching matrix. Then when  $\mathbf{C} = \mathbf{A}\mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T\mathbf{A}\mathbf{S}$ , the corresponding low-rank SPSD approximation satisfies*

$$\mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T = \mathbf{A}^{1/2}\mathbf{P}_{\mathbf{A}^{1/2}\mathbf{S}}\mathbf{A}^{1/2}.$$

##### 4.1.1 SPECTRAL NORM BOUNDS

We start with a bound on the spectral norm of the residual error. Although this result is trivial to prove given prior work, it highlights several properties that we use in the analysis of our subsequent results.

**Theorem 2** *Let  $\mathbf{A}$  be an  $n \times n$  SPSD matrix with eigenvalue decomposition partitioned as in Equation (1),  $\mathbf{S}$  be a sketching matrix of size  $n \times \ell$ ,  $q$  be a positive integer, and  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$  be as defined in Equation (3). Then when  $\mathbf{C} = \mathbf{A}^q\mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T\mathbf{A}^{2q-1}\mathbf{S}$ , the corresponding low-rank SPSD approximation satisfies*

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 \leq \|\mathbf{\Sigma}_2\|_2 + \left\| \mathbf{\Sigma}_2^{q-1/2}\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger \right\|_2^{2/(2q-1)},$$

assuming  $\mathbf{\Omega}_1$  has full row rank.

**Proof** Apply Lemma 1 with the sampling matrix  $\mathbf{S}' = \mathbf{A}^{q-1}\mathbf{S}$  (where, recall,  $q \geq 1$ ) to see that

$$\mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T = \mathbf{A}^{1/2}\mathbf{P}_{\mathbf{A}^{q-1/2}\mathbf{S}}\mathbf{A}^{1/2}.$$

It follows that

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 = \left\| \mathbf{A}^{1/2} \left( \mathbf{I} - \mathbf{P}_{(\mathbf{A}^{1/2})^{2q-1}\mathbf{S}} \right) \mathbf{A}^{1/2} \right\|_2^2. \quad (10)$$

Next, recall that  $\mathbf{\Omega}_i = \mathbf{U}_i^T \mathbf{S}$  and that  $\mathbf{A}^{1/2}$  has eigenvalue decomposition  $\mathbf{A}^{1/2} = \mathbf{U} \mathbf{\Sigma}^{1/2} \mathbf{U}^T$ , where

$$\mathbf{U} = (\mathbf{U}_1 \quad \mathbf{U}_2) \quad \text{and} \quad \mathbf{\Sigma}^{1/2} = \begin{pmatrix} \mathbf{\Sigma}_1^{1/2} & \\ & \mathbf{\Sigma}_2^{1/2} \end{pmatrix}.$$

It can be shown (see Halko et al., 2011, Theorems 9.1 and 9.2) that, because  $\mathbf{\Omega}_1$  has full row rank,

$$\left\| \mathbf{A}^{1/2} \left( \mathbf{I} - \mathbf{P}_{(\mathbf{A}^{1/2})^{2q-1} \mathbf{S}} \right) \mathbf{A}^{1/2} \right\|_2^2 \leq \left( \left\| (\mathbf{\Sigma}_2^{1/2})^{2q-1} \right\|_2^2 + \left\| (\mathbf{\Sigma}_2^{1/2})^{2q-1} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_2^2 \right)^{1/(2q-1)}. \quad (11)$$

Equations (10) and (11) imply that

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_2 &\leq \left( \left\| \mathbf{\Sigma}_2^{q-1/2} \right\|_2^2 + \left\| \mathbf{\Sigma}_2^{q-1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_2^2 \right)^{1/(2q-1)} \\ &\leq \left\| \mathbf{\Sigma}_2 \right\|_2 + \left\| \mathbf{\Sigma}_2^{q-1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_2^{2/(2q-1)} \end{aligned}$$

The latter inequality follows from the fact that the  $2q - 1$  radical function is subadditive when  $q \geq 1$  and the identity  $\left\| \mathbf{\Sigma}_2^{q-1/2} \right\|_2^2 = \left\| \mathbf{\Sigma}_2 \right\|_2^{2q-1}$ . This establishes the stated bound. ■

*Remark.* The assumption that  $\mathbf{\Omega}_1$  has full row rank is very non-trivial. It is, however, satisfied by our algorithms below. See Section 4.1.4 for more details on this point.

*Remark.* The proof of Theorem 2 proceeds in two steps. The first step relates low-rank approximation of an SPSD matrix  $\mathbf{A}$  under the SPSD Sketching Model of Section 2.2 to column sketching (*e.g.*, sampling or projecting) from the square-root of  $\mathbf{A}$ . A weaker relation of this type was used by Drineas and Mahoney (2005), but the stronger form that we use here in Equation (10) was first proved in (Gittens, 2012). The second step is to use a deterministic structural result that holds for sampling/projecting from an arbitrary matrix. The structural bound of the form of Equation (11) was originally proven for  $q = 1$  by Boutsidis et al. (2009), who applied it to the Column Subset Selection Problem. The bound was subsequently improved by Halko et al. (2011), who applied it to a random projection algorithm and extended it to apply when  $q > 1$ . Although the analyses of our next two results are more complicated, they follow the same high-level two-step approach.

#### 4.1.2 FROBENIUS NORM BOUNDS

Next, we state and prove the following bound on the Frobenius norm of the residual error. The proof parallels that for the spectral norm bound, in that we divide it into two analogous parts, but the analysis is somewhat more complex.

The multiplicative eigengap  $\gamma = \lambda_{k+1}(\mathbf{A})/\lambda_k(\mathbf{A})$  that appears in the statement of this theorem predicts the effect of using the power method when constructing sketches. Specifically, the additional errors of sketches constructed using  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$  are at least a factor of  $\gamma^{q-1}$  times smaller than those constructed using  $\mathbf{C} = \mathbf{A} \mathbf{S}$ .

**Theorem 3** Let  $\mathbf{A}$  be an  $n \times n$  SPSSD matrix with eigenvalue decomposition partitioned as in Equation (1),  $\mathbf{S}$  be a sketching matrix of size  $n \times \ell$ ,  $q$  be a positive integer,  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$  be as defined in Equation (3), and define

$$\gamma = \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})}.$$

Then when  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T \mathbf{A}^{2q-1} \mathbf{S}$ , the corresponding low-rank SPSSD approximation satisfies

$$\left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\mathbb{F}} \leq \left\| \mathbf{\Sigma}_2 \right\|_{\mathbb{F}} + \gamma^{q-1} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_2 \cdot \left( \sqrt{2 \operatorname{Tr}(\mathbf{\Sigma}_2)} + \gamma^{q-1} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\mathbb{F}} \right),$$

assuming  $\mathbf{\Omega}_1$  has full row rank.

**Proof** Apply Lemma 1 with the sampling matrix  $\mathbf{S}' = \mathbf{A}^{q-1} \mathbf{S}$  to see that

$$\mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T = \mathbf{A}^{1/2} \mathbf{P}_{\mathbf{A}^{q-1/2} \mathbf{S}} \mathbf{A}^{1/2}.$$

It follows that

$$E := \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\mathbb{F}} = \left\| \mathbf{A}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{A}^{q-1/2} \mathbf{S}}) \mathbf{A}^{1/2} \right\|_{\mathbb{F}}.$$

To bound this quantity, we first use the unitary invariance of the Frobenius norm and the fact that

$$\mathbf{P}_{\mathbf{A}^{q-1/2} \mathbf{S}} = \mathbf{U} \mathbf{P}_{\mathbf{\Sigma}^{q-1/2} \mathbf{U}^T \mathbf{S}} \mathbf{U}^T$$

to obtain

$$E^2 = \left\| \mathbf{A}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{A}^{q-1/2} \mathbf{S}}) \mathbf{A}^{1/2} \right\|_{\mathbb{F}}^2 = \left\| \mathbf{\Sigma}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{\Sigma}^{q-1/2} \mathbf{U}^T \mathbf{S}}) \mathbf{\Sigma}^{1/2} \right\|_{\mathbb{F}}^2.$$

Then we take

$$\mathbf{Z} = \mathbf{\Sigma}^{q-1/2} \mathbf{U}^T \mathbf{S} \mathbf{\Omega}_1^\dagger \mathbf{\Sigma}_1^{-(q-1/2)} = \begin{pmatrix} \mathbf{I} \\ \mathbf{F} \end{pmatrix}, \quad (12)$$

where  $\mathbf{I} \in \mathbb{R}^{k \times k}$  and  $\mathbf{F} \in \mathbb{R}^{n-k \times k}$  is given by  $\mathbf{F} = \mathbf{\Sigma}_2^{q-1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \mathbf{\Sigma}_1^{-(q-1/2)}$ . The latter equality in Equation (12) holds because of our assumption that  $\mathbf{\Omega}_1$  has full row rank. Since the range of  $\mathbf{Z}$  is contained in the range of  $\mathbf{\Sigma}^{q-1/2} \mathbf{U}^T \mathbf{S}$ ,

$$E^2 \leq \left\| \mathbf{\Sigma}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{\Sigma}^{1/2} \right\|_{\mathbb{F}}^2.$$

By construction,  $\mathbf{Z}$  has full column rank, thus  $\mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1/2}$  is an orthonormal basis for the span of  $\mathbf{Z}$ , and

$$\begin{aligned} \mathbf{I} - \mathbf{P}_{\mathbf{Z}} &= \mathbf{I} - \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T = \mathbf{I} - \begin{pmatrix} \mathbf{I} \\ \mathbf{F} \end{pmatrix} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} (\mathbf{I} \quad \mathbf{F}^T) \\ &= \begin{pmatrix} \mathbf{I} - (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} & -(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \\ -\mathbf{F}(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} & \mathbf{I} - \mathbf{F}(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \end{pmatrix}. \end{aligned} \quad (13)$$

This implies that

$$\begin{aligned}
 E^2 &\leq \left\| \Sigma^{1/2} \begin{pmatrix} \mathbf{I} - (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} & -(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \\ -\mathbf{F}(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} & \mathbf{I} - \mathbf{F}(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \end{pmatrix} \Sigma^{1/2} \right\|_{\mathbf{F}}^2 \\
 &= \left\| \Sigma_1^{1/2} (\mathbf{I} - (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1}) \Sigma_1^{1/2} \right\|_{\mathbf{F}}^2 + 2 \left\| \Sigma_1^{1/2} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \Sigma_2^{1/2} \right\|_{\mathbf{F}}^2 \\
 &\quad + \left\| \Sigma_2^{1/2} (\mathbf{I} - \mathbf{F}(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T) \Sigma_2^{1/2} \right\|_{\mathbf{F}}^2 \\
 &:= T_1 + T_2 + T_3.
 \end{aligned} \tag{14}$$

Next, we provide bounds for  $T_1$ ,  $T_2$ , and  $T_3$ . Using the fact that  $\mathbf{0} \preceq \mathbf{I} - \mathbf{F}(\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \preceq \mathbf{I}$ , we can bound  $T_3$  with

$$T_3 \leq \|\Sigma_2\|_{\mathbf{F}}^2.$$

Likewise, the fact that  $\mathbf{I} - (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \preceq \mathbf{F}^T \mathbf{F}$  (easily seen with an SVD) implies that we can bound  $T_1$  as

$$\begin{aligned}
 T_1 &\leq \left\| \Sigma_1^{1/2} \mathbf{F}^T \mathbf{F} \Sigma_1^{1/2} \right\|_{\mathbf{F}}^2 \leq \left\| \mathbf{F} \Sigma_1^{1/2} \right\|_2^2 \left\| \mathbf{F} \Sigma_1^{1/2} \right\|_{\mathbf{F}}^2 \\
 &= \left\| \Sigma_2^{q-1/2} \Omega_2 \Omega_1^\dagger \Sigma_1^{-(q-1)} \right\|_2^2 \left\| \Sigma_2^{q-1/2} \Omega_2 \Omega_1^\dagger \Sigma_1^{-(q-1)} \right\|_{\mathbf{F}}^2 \\
 &\leq \left\| \Sigma_2^{q-1} \right\|_2^4 \left\| \Sigma_1^{-(q-1)} \right\|_2^4 \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_2^2 \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\mathbf{F}}^2 \\
 &= (\|\Sigma_2\|_2 \|\Sigma_1^{-1}\|_2)^{4(q-1)} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_2^2 \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\mathbf{F}}^2 \\
 &= \left( \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})} \right)^{4(q-1)} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_2^2 \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\mathbf{F}}^2.
 \end{aligned}$$

We proceed to bound  $T_2$  by using the estimate

$$T_2 \leq 2 \left\| \Sigma_1^{1/2} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \right\|_2^2 \left\| \Sigma_2^{1/2} \right\|_{\mathbf{F}}^2. \tag{15}$$

To develop the term involving a spectral norm, observe that for any SPSD matrix  $\mathbf{M}$  with eigenvalue decomposition  $\mathbf{M} = \mathbf{V} \mathbf{D} \mathbf{V}^T$ ,

$$\begin{aligned}
 (\mathbf{I} + \mathbf{M})^{-1} \mathbf{M} (\mathbf{I} + \mathbf{M})^{-1} &= (\mathbf{V} \mathbf{V}^T + \mathbf{V} \mathbf{D} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{V}^T (\mathbf{V} \mathbf{V}^T + \mathbf{V} \mathbf{D} \mathbf{V}^T)^{-1} \\
 &= \mathbf{V} (\mathbf{I} + \mathbf{D})^{-1} \mathbf{D} (\mathbf{I} + \mathbf{D})^{-1} \mathbf{V}^T \\
 &\preceq \mathbf{V} \mathbf{D} \mathbf{V}^T = \mathbf{M}.
 \end{aligned}$$

It follows that

$$\begin{aligned}
 \left\| \boldsymbol{\Sigma}_1^{1/2} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \right\|_2^2 &= \left\| \boldsymbol{\Sigma}_1^{1/2} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{F} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \boldsymbol{\Sigma}_1^{1/2} \right\|_2^2 \\
 &\leq \left\| \boldsymbol{\Sigma}_1^{1/2} \mathbf{F}^T \mathbf{F} \boldsymbol{\Sigma}_1^{1/2} \right\|_2^2 = \left\| \mathbf{F} \boldsymbol{\Sigma}_1^{1/2} \right\|_2^2 \\
 &= \left\| \boldsymbol{\Sigma}_2^{q-1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \boldsymbol{\Sigma}_1^{-(q-1)} \right\|_2^2 \\
 &\leq \left\| \boldsymbol{\Sigma}_2^{q-1} \right\|_2^2 \left\| \boldsymbol{\Sigma}_1^{-(q-1)} \right\|_2^2 \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2^2 \\
 &= \left( \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})} \right)^{2(q-1)} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2^2.
 \end{aligned}$$

Using this estimate in Equation (15), we conclude that

$$T_2 \leq 2 \left( \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})} \right)^{2(q-1)} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2^2 \left\| \boldsymbol{\Sigma}_2^{1/2} \right\|_F^2.$$

Combining our estimates for  $T_1$ ,  $T_2$ , and  $T_3$  with Equation (14) gives

$$\begin{aligned}
 E^2 &= \left\| \mathbf{A}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{A}^{q-1/2} \mathbf{S}}) \mathbf{A}^{1/2} \right\|_F^2 \leq \left\| \boldsymbol{\Sigma}_2 \right\|_F^2 \\
 &\quad + \left( \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})} \right)^{2(q-1)} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2^2 \cdot \left( 2 \left\| \boldsymbol{\Sigma}_2^{1/2} \right\|_F^2 + \left( \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})} \right)^{2(q-1)} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_F^2 \right).
 \end{aligned}$$

The claimed bound follows by identifying  $\gamma$  and applying the subadditivity of the square-root function:

$$E \leq \left\| \boldsymbol{\Sigma}_2 \right\|_F + \gamma^{q-1} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2 \cdot \left( \sqrt{2 \operatorname{Tr}(\boldsymbol{\Sigma}_2)} + \gamma^{q-1} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_F \right). \quad \blacksquare$$

*Remark.* The quality of approximation guarantee provided by Theorem 3 depends on the quantities  $\left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2$  and  $\left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_F$ ; these quantities reflect the extent to which the sketching matrix is aligned with the eigenspaces of  $\mathbf{A}$ . The dependence on  $\gamma$  captures the facts that the power method is effective only when there is spectral decay, and that larger gaps between the  $k$  and  $k+1$  eigenvalues lead to smaller errors when the power method is used.

*Remark.* To obtain a greater understanding of the additional error term in Theorem 3, assume that  $\mathbf{S}$  is a particularly effective sketching matrix, so that  $\left\| \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2 = \mathcal{O}(1)$ . Then

$$\left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2 = \mathcal{O} \left( \left\| \boldsymbol{\Sigma}_2 \right\|_2^{1/2} \right) \quad \text{and} \quad \sqrt{2 \operatorname{Tr}(\boldsymbol{\Sigma}_2)} + \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_F = \mathcal{O} \left( \left\| \boldsymbol{\Sigma}_2 \right\|_*^{1/2} \right),$$

and the theorem guarantees that the additional error is on the order of  $\sqrt{\left\| \boldsymbol{\Sigma}_2 \right\|_2 \left\| \boldsymbol{\Sigma}_2 \right\|_*}$ . This is an upper bound on the optimal Frobenius error:

$$\left\| \boldsymbol{\Sigma}_2 \right\|_F \leq \sqrt{\left\| \boldsymbol{\Sigma}_2 \right\|_2 \left\| \boldsymbol{\Sigma}_2 \right\|_*}.$$

We see, in particular, that if the residual spectrum is flat, i.e.  $\lambda_{k+1}(\mathbf{A}) = \dots = \lambda_n(\mathbf{A})$ , then equality holds and the additional error is on the scale of the optimal error.

## 4.1.3 TRACE NORM BOUNDS

Finally, we state and prove the following bound on the trace norm of the residual error. The proof method is analogous to that for the spectral and Frobenius norm bounds.

As in the case of the Frobenius norm error, we see that the multiplicative eigengap  $\gamma = \lambda_{k+1}(\mathbf{A})/\lambda_k(\mathbf{A})$  predicts the effect of using the power method when constructing sketches.

**Theorem 4** *Let  $\mathbf{A}$  be an  $n \times n$  SPSD matrix with eigenvalue decomposition partitioned as in Equation (1),  $\mathbf{S}$  be a sketching matrix of size  $n \times \ell$ ,  $q$  be a positive integer,  $\mathbf{\Omega}_1$  and  $\mathbf{\Omega}_2$  be as defined in Equation (3), and define*

$$\gamma = \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})}.$$

*Then when  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T \mathbf{A}^{2q-1} \mathbf{S}$ , the corresponding low-rank SPSD approximation satisfies*

$$\left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\star} \leq \text{Tr}(\mathbf{\Sigma}_2) + \gamma^{2(q-1)} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\mathbf{F}}^2,$$

*assuming  $\mathbf{\Omega}_1$  has full row rank.*

**Proof** Since  $\mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T = \mathbf{A}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{A}^{q-1/2} \mathbf{S}}) \mathbf{A}^{1/2} \succeq \mathbf{0}$ , its trace norm simplifies to its trace. Thus

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\star} &= \text{Tr}(\mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T) = \text{Tr}(\mathbf{\Sigma}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{\Sigma}^{q-1/2} \mathbf{S}}) \mathbf{\Sigma}^{1/2}) \\ &\leq \text{Tr}(\mathbf{\Sigma}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{\Sigma}^{1/2}), \end{aligned}$$

where  $\mathbf{Z} = \begin{pmatrix} \mathbf{I} \\ \mathbf{F} \end{pmatrix}$  is defined in Equation (12). The expression for  $\mathbf{I} - \mathbf{P}_{\mathbf{Z}}$  given in Equation (13) implies that

$$\begin{aligned} &\text{Tr}(\mathbf{\Sigma}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{\Sigma}^{1/2}) \\ &= \text{Tr}(\mathbf{\Sigma}_1^{1/2} (\mathbf{I} - (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1}) \mathbf{\Sigma}_1^{1/2}) + \text{Tr}(\mathbf{\Sigma}_2^{1/2} (\mathbf{I} - \mathbf{F} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T) \mathbf{\Sigma}_2^{1/2}). \end{aligned}$$

Recall the estimate  $\mathbf{I} - (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \preceq \mathbf{F}^T \mathbf{F}$  and the basic estimate  $\mathbf{I} - \mathbf{F} (\mathbf{I} + \mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \preceq \mathbf{I}$ . Together these imply that

$$\begin{aligned} \text{Tr}(\mathbf{\Sigma}^{1/2} (\mathbf{I} - \mathbf{P}_{\mathbf{Z}}) \mathbf{\Sigma}^{1/2}) &\leq \text{Tr}(\mathbf{\Sigma}_1^{1/2} \mathbf{F}^T \mathbf{F} \mathbf{\Sigma}_1^{1/2}) + \text{Tr}(\mathbf{\Sigma}_2) \\ &= \text{Tr}(\mathbf{\Sigma}_2) + \left\| \mathbf{\Sigma}_2^{q-1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \mathbf{\Sigma}_1^{-(q-1)} \right\|_{\mathbf{F}}^2 \\ &\leq \text{Tr}(\mathbf{\Sigma}_2) + \left\| \mathbf{\Sigma}_2^{q-1} \right\|_2^2 \left\| \mathbf{\Sigma}_1^{-(q-1)} \right\|_2^2 \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\mathbf{F}}^2 \\ &= \text{Tr}(\mathbf{\Sigma}_2) + \gamma^{2(q-1)} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\mathbf{F}}^2. \end{aligned}$$



The first equality follows from substituting the definition of  $\mathbf{F}$  and identifying the squared Frobenius norm. The last equality follows from identifying  $\gamma$ . We have established the claimed bound. ■

*Remark.* Since the identity  $\|\mathbf{X}\|_F^2 = \|\mathbf{X}\mathbf{X}^T\|_*$  holds for any matrix  $\mathbf{X}$ , the squared Frobenius norm term present in the deterministic error bound for the trace norm error is on the scale of  $\|\boldsymbol{\Sigma}_2\|_*$  when  $\|\boldsymbol{\Omega}_2\boldsymbol{\Omega}_1^\dagger\|_2$  is  $O(1)$ .

#### 4.1.4 ADDITIONAL REMARKS ON OUR DETERMINISTIC STRUCTURAL RESULTS

Before applying these deterministic structural results in particular randomized algorithmic settings, we pause to make several additional remarks about these three theorems.

First, for some randomized sampling schemes, it may be difficult to obtain a sharp bound on  $\|\boldsymbol{\Omega}_2\boldsymbol{\Omega}_1^\dagger\|_\xi$  for  $\xi = 2, F$ . In these situations, the bounds on the excess error supplied by Theorems 2, 3, and 4 may be quite pessimistic. On the other hand, since  $\mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T = \mathbf{A}^{1/2}(\mathbf{I} - \mathbf{P}_{(\mathbf{A}^{1/2})^{2q-1}\mathbf{S}})\mathbf{A}^{1/2}$ , it follows that  $\mathbf{0} \preceq \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \preceq \mathbf{A}$ . This implies that the errors of *any* approximation generated using the SPSD Sketching Model, deterministic or randomized, satisfy at least the crude bound  $\|\mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T\|_\xi \leq \|\mathbf{A}\|_\xi$ .

Second, we emphasize that these theorems are deterministic structural results that bound the additional error (beyond that of the optimal rank- $k$  approximation) of low-rank approximations which follow our SPSD sketching model. That is, there is no randomness in their statement or analysis. In particular, these bounds hold for deterministic as well as randomized sketching matrices  $\mathbf{S}$ . In the latter case, the randomness enters only through  $\mathbf{S}$ , and one needs to show that the condition that  $\boldsymbol{\Omega}_1$  has full row rank is satisfied with high probability; conditioned on this, the quality of the bound is determined by terms that depend on how the sketching matrix interacts with the subspace structure of the matrix  $\mathbf{A}$ .

In particular, we remind the reader that (although it is beyond the scope of this paper to explore this point in detail) these deterministic structural results could be used to check, in an *a posteriori* manner, the quality of a sketching method for which one cannot establish an *a priori* bound.

Third, we also emphasize that the assumption that  $\boldsymbol{\Omega}_1$  has full row rank (equivalently, that  $\tan(\mathbf{S}, \mathbf{U}_1) < \infty$ ) is very non-trivial; and that it is false, in worst-case at least and for non-trivial parameter values, for common sketching methods such as uniform sampling. To see that some version of leverage-based sampling is needed to ensure this condition, recall that  $\mathbf{U}_1^T\mathbf{U}_1 = \mathbf{I}$  and thus that  $\boldsymbol{\Omega}_1\boldsymbol{\Omega}_1^T = \mathbf{U}_1^T\mathbf{S}\mathbf{S}^T\mathbf{U}_1$  can be viewed as approximating  $\mathbf{I}$  with a small number of rank-1 components of  $\mathbf{U}_1^T\mathbf{U}_1$ . The condition that  $\boldsymbol{\Omega}_1$  has full row rank is equivalent to  $\|\mathbf{U}_1^T\mathbf{U}_1 - \mathbf{U}_1^T\mathbf{S}\mathbf{S}^T\mathbf{U}_1\|_2 < 1$ . Work on approximating the product of matrices by random sampling shows that to obtain non-trivial bounds one must sample with respect to the norm of the rank-1 components (Drineas et al., 2006), which here (since we are approximating the product of two orthogonal matrices) equal the statistical leverage scores. From this perspective, random projections satisfy this condition since (informally) they rotate to a random basis where the leverage scores of the rotated matrix are approximately uniform and thus where uniform sampling is appropriate (Drineas et al., 2010; Mahoney, 2011).

Finally, as observed recently in Bach (2013), methods that use knowledge of a matrix square root  $\Phi$  (*i.e.*, a  $\Phi$  such that  $\mathbf{A} = \Phi\Phi^T$ ) typically lead to  $\Omega(n^2)$  complexity. An important feature of our approach is that we only use the matrix square root implicitly—that is, inside the analysis, and not in the statement of the algorithm—and thus we do *not* incur any such cost.

## 4.2 Stochastic Error Bounds for Low-rank SPSD Approximation

In this section, we apply the three theorems from Section 4.1 to bound the reconstruction errors for several random sampling and random projection methods that conform to our SPSD Sketching Model. In particular, we consider two variants of random sampling and two variants of random projections: sampling columns according to an importance sampling distribution that depends on the statistical leverage scores (in Section 4.2.1); randomly projecting by using subsampled randomized Fourier transformations (in Section 4.2.2); randomly projecting by uniformly sampling from Gaussian mixtures of the columns (in Section 4.2.3); and, finally, sampling columns uniformly at random (in Section 4.2.4).

The results are presented for the general case of SPSD sketches constructed using the power method, *i.e.*, sketches constructed using  $\mathbf{C} = \mathbf{A}^q\mathbf{S}$  for a positive integer  $q$ . The additive errors of these sketches decrease proportionally to the number of iterations  $q$ , where the constant of proportionality is given by the multiplicative eigengap  $\gamma = \lambda_{k+1}(\mathbf{A})/\lambda_k(\mathbf{A})$ . Accordingly, the bounds involve the terms  $\gamma^{q-1}$  and  $\gamma^{2(q-1)}$ . The bounds simplify considerably when  $q = 1$  (*i.e.*, when there are no additional iterations) or  $\gamma = 1$  (*i.e.*, when there is no eigengap). In either of these cases, the terms  $\gamma^{q-1}$  and  $\gamma^{2(q-1)}$  all become the constant 1.

Before establishing these results, we pause here to provide a brief review of running time issues, some of which were addressed empirically in Section 3. The computational bottleneck for random sampling algorithms (except for uniform sampling that we address in Section 4.2.4, which is trivial to implement) is often the exact or approximate computation of the importance sampling distribution with respect to which one samples; and the computational bottleneck for random projection methods is often the implementation of the random projection. For example, if the sketching matrix  $\mathbf{S}$  is a random projection constructed as an  $n \times \ell$  matrix of i.i.d. Gaussian random variables, as we use in Section 4.2.3, then the running time of dense data in RAM is not substantially faster than computing  $\mathbf{U}_1$ , while the running time can be much faster for certain sparse matrices or for computation in parallel or distributed environments. Alternately, if the sketching matrix  $\mathbf{S}$  is a Fourier-based projection, as we use in Section 4.2.2, then the running time for data stored in RAM is typically  $O(n^2 \ln k)$ , as opposed to the  $O(n^2 k)$  time that would be needed to compute  $\mathbf{U}_1$ . These running times depend sensitively on the size of the data and the model of data access; see Mahoney (2011); Halko et al. (2011) for detailed discussions of these issues.

In particular, for random sampling algorithms that use a leverage-based importance sampling distribution, as we use in Section 4.2.1, it is often said that the running time is no faster than that of computing  $\mathbf{U}_1$ . (This  $O(n^2 k)$  running time claim is simply the running time of the naïve algorithm that computes  $\mathbf{U}_1$  “exactly,” *e.g.*, with a variant of the QR decomposition, and then reads off the Euclidean norms of the rows.) However, the randomized algorithm of Drineas et al. (2012) that computes relative-error approximations to *all* of the statistical leverage in a time that is qualitatively faster—in worst-case theory

and, by using existing high-quality randomized numerical code (Avron et al., 2010; Meng et al., 2014; Halko et al., 2011), in practice—gets around this bottleneck, as was shown in Section 3. The computational bottleneck for the algorithms of Drineas et al. (2012) is that of applying a random projection, and thus the running time for leverage-based Nyström extension is that of applying a (“fast” Fourier-based or “slow” Gaussian-based, as appropriate) random projection to  $\mathbf{A}$  (Drineas et al., 2012). See Section 3 or (Avron et al., 2010; Meng et al., 2014; Halko et al., 2011) for additional details.

#### 4.2.1 SAMPLING WITH LEVERAGE-BASED IMPORTANCE SAMPLING PROBABILITIES

Here, the columns of  $\mathbf{A}$  are sampled with replacement according to a nonuniform probability distribution determined by the (exact or approximate) statistical leverage scores of  $\mathbf{A}$  relative to the best rank- $k$  approximation to  $\mathbf{A}$ , which in turn depend on nonuniformity properties of the top  $k$ -dimensional eigenspace of  $\mathbf{A}$ . To add flexibility (*e.g.*, in case the scores are computed only approximately with the fast algorithm of Drineas et al. (2012)), we formulate the following lemma in terms of any probability distribution that is  $\beta$ -close to the leverage score distribution. In particular, consider any probability distribution satisfying

$$p_j \geq \frac{\beta}{k} \|(\mathbf{U}_1)_j\|_2^2 \quad \text{and} \quad \sum_{j=1}^n p_j = 1,$$

where  $\beta \in (0, 1]$ . Given these ( $\beta$ -approximate) leverage-based probabilities, the sketching matrix is  $\mathbf{S} = \mathbf{R}\mathbf{D}$  where  $\mathbf{R} \in \mathbb{R}^{n \times \ell}$  is a column selection matrix that samples columns of  $\mathbf{A}$  from the given distribution—*i.e.*,  $\mathbf{R}_{ij} = 1$  iff the  $i$ th column of  $\mathbf{A}$  is the  $j$ th column selected—and  $\mathbf{D}$  is a diagonal rescaling matrix satisfying  $\mathbf{D}_{jj} = \frac{1}{\sqrt{\ell p_i}}$  iff  $\mathbf{R}_{ij} = 1$ . For this case, we can prove the following.

**Lemma 5** *Let  $\mathbf{A}$  be an  $n \times n$  SPSD matrix,  $q$  be a positive integer, and  $\mathbf{S}$  be a sampling matrix of size  $n \times \ell$  corresponding to a leverage-based probability distribution derived from the top  $k$ -dimensional eigenspace of  $\mathbf{A}$ , satisfying*

$$p_j \geq \frac{\beta}{k} \|(\mathbf{U}_1)_j\|_2^2 \quad \text{and} \quad \sum_{j=1}^n p_j = 1$$

for some  $\beta \in (0, 1]$ . Fix a failure probability  $\delta \in (0, 1]$  and approximation factor  $\epsilon \in (0, 1]$ , and let

$$\gamma = \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})}.$$

If  $\ell \geq 3200(\beta\epsilon^2)^{-1}k \ln(4k/(\beta\delta))$ , then, when  $\mathbf{C} = \mathbf{A}^q\mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T\mathbf{A}^{2q-1}\mathbf{S}$ , the corresponding low-rank SPSD approximation satisfies

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 \leq \|\mathbf{A} - \mathbf{A}_k\|_2 + (\epsilon^2 \|(\mathbf{A} - \mathbf{A}_k)^{2q-1}\|_\star)^{1/(2q-1)}, \quad (16)$$

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_F \leq \|\mathbf{A} - \mathbf{A}_k\|_F + \left( \sqrt{2}\epsilon\gamma^{q-1} + \epsilon^2\gamma^{2(q-1)} \right) \|\mathbf{A} - \mathbf{A}_k\|_\star, \quad \text{and} \quad (17)$$

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_\star \leq (1 + \gamma^{2(q-1)}\epsilon^2) \|\mathbf{A} - \mathbf{A}_k\|_\star, \quad (18)$$

simultaneously with probability at least  $1 - 6\delta - 0.6$ .

**Proof** In (Mackey et al., 2011a, proof of Proposition 22) it is shown that if  $\ell$  satisfies the given bound and the samples are drawn from an approximate subspace probability distribution, then for any SPSD diagonal matrix  $\mathbf{D}$ ,

$$\left\| \mathbf{D} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_{\text{F}} \leq \epsilon \|\mathbf{D}\|_{\text{F}}$$

with probability at least  $1 - 2\delta - 0.2$ . Thus, the estimates

$$\left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_{\text{F}} \leq \epsilon \left\| \boldsymbol{\Sigma}_2^{1/2} \right\|_{\text{F}} = \epsilon \sqrt{\text{Tr}(\boldsymbol{\Sigma}_2)} = \epsilon \sqrt{\|\mathbf{A} - \mathbf{A}_k\|_{\star}},$$

and

$$\begin{aligned} \left( \left\| \boldsymbol{\Sigma}_2^{q-1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2 \right)^{2/(2q-1)} &\leq \left( \left\| \boldsymbol{\Sigma}_2^{p-1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_{\text{F}} \right)^{2/(2q-1)} \\ &\leq \left( \epsilon^2 \left\| \boldsymbol{\Sigma}_2^{q-1/2} \right\|_{\text{F}}^2 \right)^{1/(2q-1)} \\ &= \left( \epsilon^2 \text{Tr}(\boldsymbol{\Sigma}_2^{2q-1}) \right)^{1/(2q-1)} \\ &= \left( \epsilon^2 \|\mathbf{A} - \mathbf{A}_k\|_{\star}^{2q-1} \right)^{1/(2q-1)} \end{aligned}$$

each hold, individually, with probability at least  $1 - 2\delta - 0.2$ . In particular, taking  $q = 1$ , we see that

$$\left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2 \leq \epsilon \sqrt{\|\mathbf{A} - \mathbf{A}_k\|_{\star}}$$

with the same probability.

These three estimates used in Theorems 2, 3, and 4 yield the bounds given in the statement of the theorem.  $\blacksquare$

*Remark.* The additive scale factors for the spectral and Frobenius norm bounds are much improved relative to the prior results of Drineas and Mahoney (2005). At root, this is since the leverage score importance sampling probabilities highlight structural properties of the data (*e.g.*, how to satisfy the condition in Theorems 2, 3, and 4 that  $\boldsymbol{\Omega}_1$  has full row rank) in a more refined way than the importance sampling probabilities of Drineas and Mahoney (2005).

*Remark.* These improvements come at additional computational expense, but we remind the reader that leverage-based sampling probabilities of the form used by Lemma 5 can be computed faster than the time needed to compute the basis  $\mathbf{U}_1$  (Drineas et al., 2012). The computational bottleneck of the algorithm of Drineas et al. (2012) is the time required to perform a random projection on the input matrix.

*Remark.* Not surprisingly, constant factors such as 3200 (as well as other similarly large factors below) and a failure probability bounded away from zero are artifacts of the analysis; the empirical behavior of this sampling method is much better. This has been observed previously (Drineas et al., 2008; Mahoney and Drineas, 2009).

## 4.2.2 RANDOM PROJECTIONS WITH SUBSAMPLED RANDOMIZED FOURIER TRANSFORMS

Here, the columns of  $\mathbf{A}$  are randomly mixed using a unitary matrix before the columns are sampled. In particular,  $\mathbf{S} = \sqrt{\frac{n}{\ell}} \mathbf{D} \mathbf{T} \mathbf{R}$ , where  $\mathbf{D}$  is a diagonal matrix of Rademacher random variables,  $\mathbf{T}$  is a highly incoherent unitary matrix, and  $\mathbf{R}$  restricts to  $\ell$  columns. For concreteness, and because it has an associated fast transform, we consider the case where  $\mathbf{T}$  is the normalized Fourier transform of size  $n \times n$ . For this case, we can prove the following.

**Lemma 6** *Let  $\mathbf{A}$  be an  $n \times n$  SPSD matrix,  $q$  be a positive integer, and  $\mathbf{S} = \sqrt{\frac{n}{\ell}} \mathbf{D} \mathbf{F} \mathbf{R}$  be a sampling matrix of size  $n \times \ell$ , where  $\mathbf{D}$  is a diagonal matrix of Rademacher random variables,  $\mathbf{F}$  is a normalized Fourier matrix of size  $n \times n$ , and  $\mathbf{R}$  restricts to  $\ell$  columns. Fix a failure probability  $\delta \in (0, 1)$ , approximation factor  $\epsilon \in (0, 1)$ , and assume that  $k \geq 4$ . Define*

$$\gamma = \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})}.$$

*If  $\ell \geq 24\epsilon^{-1}[\sqrt{k} + \sqrt{8 \ln(8n/\delta)}]^2 \ln(8k/\delta)$ , then, when  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T \mathbf{A}^{2q-1} \mathbf{S}$ , the corresponding low-rank SPSD approximation satisfies*

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_2 &\leq \left[ 1 + \left( \frac{1}{1 - \sqrt{\epsilon}} \cdot \left( 5 + \frac{16 \ln(n/\delta)^2}{\ell} \right) \right)^{1/(2q-1)} \right] \cdot \|\mathbf{A} - \mathbf{A}_k\|_2 \\ &\quad + \left( \frac{2 \ln(n/\delta)}{(1 - \sqrt{\epsilon}) \ell} \right)^{1/(2q-1)} \left\| (\mathbf{A} - \mathbf{A}_k)^{2q-1} \right\|_\star^{1/(2q-1)}, \quad (19) \\ \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\mathbb{F}} &\leq \|\mathbf{A} - \mathbf{A}_k\|_{\mathbb{F}} + (7\gamma^{q-1} \sqrt{\epsilon} + 22\gamma^{2q-2} \epsilon) \|\mathbf{A} - \mathbf{A}_k\|_\star, \text{ and} \\ \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_\star &\leq (1 + 22\epsilon\gamma^{2(q-1)}) \|\mathbf{A} - \mathbf{A}_k\|_\star \end{aligned}$$

*simultaneously with probability at least  $1 - 2\delta$ .*

**Proof** Let  $\mathbf{M} = \Sigma_2^{q-1/2} \Omega_2 \Omega_1^\dagger$  denote the matrix referenced in Theorems 2 and 4. In (Boutsidis and Gittens, 2013, proof of Theorem 4), it is shown that for the stated choice of  $\mathbf{S}$  and number of samples  $\ell$ ,

$$\begin{aligned} \|\mathbf{M}\|_2^2 &\leq \frac{1}{1 - \sqrt{\epsilon}} \cdot \left( 5 \|\Sigma_2\|_2^{2q-1} + \frac{\ln(n/\delta)}{\ell} \left( \left\| \Sigma_2^{q-1/2} \right\|_{\mathbb{F}} + \sqrt{8 \ln(n/\delta)} \left\| \Sigma_2^{q-1/2} \right\|_2 \right)^2 \right) \\ &= \frac{1}{1 - \sqrt{\epsilon}} \cdot \left( 5 \|\Sigma\|_2^{2q-1} + \frac{\ln(n/\delta)}{\ell} \left( \left\| \Sigma_2^{2q-1} \right\|_\star^{1/2} + \sqrt{8 \ln(n/\delta)} \|\Sigma_2\|_2^{q-1/2} \right)^2 \right) \\ &\leq \frac{1}{1 - \sqrt{\epsilon}} \cdot \left( \left( 5 + \frac{16 \ln(n/\delta)^2}{\ell} \right) \|\Sigma_2\|_2^{2q-1} + \frac{2 \ln(n/\delta)}{\ell} \left\| \Sigma_2^{2q-1} \right\|_\star \right) \end{aligned}$$

and

$$\|\mathbf{M}\|_{\mathbb{F}} \leq \sqrt{22\epsilon} \left\| \Sigma_2^{1/2} \right\|_{\mathbb{F}} = \sqrt{22\epsilon} \|\Sigma_2\|_\star$$

each hold, individually, with probability at least  $1 - \delta$ . These estimates used in Theorems 2 and 4 yield the stated bounds for the spectral and trace norm errors.

The Frobenius norm bound follows from the same estimates and a simplification of the bound stated in Theorem 3:

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_{\mathbb{F}} &\leq \|\boldsymbol{\Sigma}_2\|_{\mathbb{F}} + \gamma^{q-1} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_2 \left( \sqrt{2 \operatorname{Tr}(\boldsymbol{\Sigma}_2)} + \gamma^{q-1} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_{\mathbb{F}} \right) \\ &\leq \|\boldsymbol{\Sigma}_2\|_{\mathbb{F}} + \gamma^{q-1} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_{\mathbb{F}} \sqrt{2 \operatorname{Tr}(\boldsymbol{\Sigma}_2)} + \gamma^{2(q-1)} \left\| \boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger \right\|_{\mathbb{F}}^2 \\ &\leq \|\boldsymbol{\Sigma}_2\|_{\mathbb{F}} + \left( \gamma^{q-1} \sqrt{44\epsilon} + 22\gamma^{2q-2}\epsilon \right) \|\boldsymbol{\Sigma}_2\|_{\star}. \end{aligned}$$

We note that a direct application of Theorem 3 gives a potentially tighter, but more unwieldy, bound.  $\blacksquare$

*Remark.* Suppressing the dependence on  $\delta$  and  $\epsilon$ , the spectral norm bound ensures that when  $q = 1$ ,  $k = \Omega(\ln n)$  and  $\ell = \Omega(k \ln k)$ , then

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 = \mathcal{O} \left( \frac{\ln n}{\ln k} \|\mathbf{A} - \mathbf{A}_k\|_2 + \frac{1}{\ln k} \|\mathbf{A} - \mathbf{A}_k\|_{\star} \right).$$

This should be compared to the guarantee established in Lemma 7 below for Gaussian-based SPSD sketches constructed using the same number of measurements:

$$\left\| \mathbf{A} - \mathbf{C}\mathbf{W}^\dagger\mathbf{C}^T \right\|_2 = \mathcal{O} \left( \|\mathbf{A} - \mathbf{A}_k\|_2 + \frac{1}{k \ln k} \|\mathbf{A} - \mathbf{A}_k\|_{\star} \right).$$

Lemma 6 guarantees that errors on this order can be achieved if one increases the number of samples by a logarithm factor in the dimension: specifically, such a bound is achieved when  $k = \Omega(\ln n)$  and  $\ell = \Omega(k \ln k \ln n)$ . The difference between the number of samples necessary for Fourier-based sketches and Gaussian-based sketches is reflective of the differing natures of the random projections: the geometry of any  $k$ -dimensional subspace is preserved under projection onto the span of  $\ell = \mathcal{O}(k)$  Gaussian random vectors (Halko et al., 2011), but the sharpest analysis available suggests that to preserve the geometry of such a subspace under projection onto the span of  $\ell$  SRFT vectors,  $\ell$  must satisfy  $\ell = \Omega(\max\{k, \ln n\} \ln k)$  (Tropp, 2011). We note, however, that in practice the Fourier-based and Gaussian-based SPSD sketches have similar reconstruction errors.

*Remark.* The structure of the Frobenius and trace norm bounds for the Fourier-based projection are identical to the structure of the corresponding bounds from Lemma 5 for leverage-based sampling (and the bounds could be made identical with appropriate choice of parameters). This is not surprising since (informally) Fourier-based (and other) random projections rotate to a random basis where the leverage scores are approximately uniform and thus where uniform sampling is appropriate (Mahoney, 2011). The disparity of the spectral norm bounds suggests that leverage-based SPSD sketches should be expected to be more accurate in the spectral norm than Fourier-based sketches; the empirical results of Section 3.4 support this interpretation. The running times of the Fourier-based and the leverage-based algorithms are the same, to leading order, if the algorithm of Drineas et al. (2012) (which uses the same transform  $\mathbf{S} = \sqrt{\frac{n}{\ell}} \mathbf{D}\mathbf{T}\mathbf{R}$ ) is used to approximate the leverage scores.

## 4.2.3 RANDOM PROJECTIONS WITH I.I.D. GAUSSIAN RANDOM MATRICES

Here, the columns of  $\mathbf{A}$  are randomly mixed using Gaussian random variables before sampling. Thus, the entries of the sampling matrix  $\mathbf{S} \in \mathbb{R}^{n \times \ell}$  are i.i.d. standard Gaussian random variables.

**Lemma 7** *Let  $\mathbf{A}$  be an  $n \times n$  SPSD matrix,  $q$  be a positive integer,  $\mathbf{S} \in \mathbb{R}^{n \times \ell}$  be a matrix of i.i.d. standard Gaussians, and define*

$$\gamma = \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})}.$$

*If  $\ell \geq 2\epsilon^{-2}k \ln k$  where  $\epsilon \in (0, 1)$  and  $k > 4$ , then, when  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T \mathbf{A}^{2q-1} \mathbf{S}$ , the corresponding low-rank SPSD approximation satisfies*

$$\begin{aligned} \|\mathbf{A} - \mathbf{C}\mathbf{W}^\dagger \mathbf{C}^T\|_2 &\leq \left(1 + \left(89 \frac{\epsilon^2}{\ln k} + 874 \frac{\epsilon^2}{k}\right)^{1/(2q-1)}\right) \|\mathbf{A} - \mathbf{A}_k\|_2 \\ &\quad + \left(219 \frac{\epsilon^2}{k \ln k}\right)^{1/(2q-1)} \cdot \|\mathbf{A} - \mathbf{A}_k\|_\star, \\ \|\mathbf{A} - \mathbf{C}\mathbf{W}^\dagger \mathbf{C}^T\|_F &\leq \|\mathbf{A} - \mathbf{A}_k\|_F + \left[\gamma^{q-1} \epsilon \left(\frac{42}{\sqrt{k}} + \frac{14}{\sqrt{\ln k}}\right) \right. \\ &\quad \left. + \gamma^{2q-2} \epsilon^2 \left(\frac{45}{\ln k} + \frac{140}{\sqrt{k \ln k}} + \frac{219}{k \sqrt{\ln k}}\right)\right] \sqrt{\|\mathbf{A} - \mathbf{A}_k\|_2 \|\mathbf{A} - \mathbf{A}_k\|_\star} \\ &\quad + \left(21 \gamma^{q-1} \frac{\epsilon}{\sqrt{k \ln k}} + 70 \gamma^{2q-2} \frac{\epsilon^2}{\sqrt{k \ln k}}\right) \|\mathbf{A} - \mathbf{A}_k\|_\star \\ &\quad + \gamma^{2q-2} \epsilon^2 \left(\frac{140}{\sqrt{k \ln k}} + \frac{437}{k}\right) \|\mathbf{A} - \mathbf{A}_k\|_2, \text{ and} \\ \|\mathbf{A} - \mathbf{C}\mathbf{W}^\dagger \mathbf{C}^T\|_\star &\leq \left(1 + 45 \frac{\gamma^{2q-2} \epsilon^2}{\ln k}\right) \|\mathbf{A} - \mathbf{A}_k\|_\star + 437 \frac{\gamma^{2q-2} \epsilon^2}{k} \|\mathbf{A} - \mathbf{A}_k\|_2 \end{aligned}$$

*simultaneously with probability at least  $1 - 2k^{-1} - 4k^{-k/\epsilon^2}$ .*

*Remark.* The way we have parameterized these bounds for Gaussian-based projections makes explicit the dependence on various parameters, but hides the structural simplicity of these bounds. In particular, note that the Frobenius norm approximation error is upper bounded by a term that depends on the Frobenius norm error of the optimal low-rank approximant and a term that depends on the trace norm error of the optimal low-rank approximant; and that, similarly, the trace norm approximation error is upper bounded by a multiplicative factor that can be set to  $1 + \epsilon$  with an appropriate choice of parameters.

**Proof** As before, this result is established by bounding the quantities involved in Theorems 2, 3, and 4. The following deviation bounds, established in (Halko et al., 2011, Section 10), are useful in that regard: if  $\mathbf{D}$  is a diagonal matrix,  $\ell = k + p$  with  $p > 4$  and  $u, t \geq 1$ ,

then with our choice of  $\mathbf{S}$ ,

$$\begin{aligned} \mathbb{P} \left\{ \left\| \mathbf{D}\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger \right\|_2 > \|\mathbf{D}\|_2 \left( \sqrt{\frac{3k}{p+1}} \cdot t + \frac{e\sqrt{\ell}}{p+1} \cdot tu \right) + \|\mathbf{D}\|_{\mathbb{F}} \frac{e\sqrt{\ell}}{p+1} \cdot t \right\} &\leq 2t^{-p} + e^{-u^2/2}, \text{ and} \\ \mathbb{P} \left\{ \left\| \mathbf{D}\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger \right\|_{\mathbb{F}} > \|\mathbf{D}\|_{\mathbb{F}} \sqrt{\frac{3k}{p+1}} \cdot t + \|\mathbf{D}\|_2 \frac{e\sqrt{\ell}}{p+1} \cdot tu \right\} &\leq 2t^{-p} + e^{-u^2/2}. \end{aligned} \quad (20)$$

Write  $\ell = k + p$ . Since  $\ell \geq 2\epsilon^{-2}k \ln k$ , we have that  $p \geq \epsilon^{-2}k \ln k$ . Accordingly, the following estimates hold:

$$\begin{aligned} \sqrt{\frac{3k}{p+1}} &\leq \sqrt{\frac{3k}{p}} \leq \sqrt{\frac{3}{\ln k}} \epsilon \\ \frac{\sqrt{\ell}}{p+1} &\leq \frac{\sqrt{k+p}}{p} \leq \sqrt{\frac{\epsilon^4}{k \ln^2 k} + \frac{\epsilon^2}{k \ln k}} < \sqrt{\frac{2}{k \ln k}} \epsilon. \end{aligned}$$

Use these estimates and take  $t = e$  and  $u = \sqrt{2 \ln k}$  in (20) to obtain that

$$\begin{aligned} \left\| \mathbf{\Sigma}_2^{q-1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_2^2 &\leq \left[ \epsilon \left( e \sqrt{\frac{3}{\ln k}} + 2e^2 \sqrt{\frac{1}{k}} \right) \cdot \left\| \mathbf{\Sigma}_2^{q-1/2} \right\|_2 + \epsilon e^2 \sqrt{\frac{2}{k \ln k}} \cdot \left\| \mathbf{\Sigma}_2^{q-1/2} \right\|_{\mathbb{F}} \right]^2 \\ &\leq 2\epsilon^2 \left( e \sqrt{\frac{3}{\ln k}} + 2e^2 \sqrt{\frac{1}{k}} \right)^2 \cdot \left\| \mathbf{\Sigma}_2 \right\|_2^{2q-1} + \frac{4\epsilon^2 e^4}{k \ln k} \cdot \left\| \mathbf{\Sigma}_2^{q-1/2} \right\|_{\mathbb{F}}^2 \\ &\leq \left( \frac{12e^2}{\ln k} + \frac{16e^4}{k} \right) \epsilon^2 \cdot \left\| \mathbf{\Sigma}_2 \right\|_2^{2q-1} + \frac{4\epsilon^2 e^4}{k \ln k} \cdot \left\| \mathbf{\Sigma}_2^{2q-1} \right\|_{\star} \end{aligned}$$

with probability at least  $1 - k^{-1} - 2k^{-k/\epsilon^2}$  and

$$\begin{aligned} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\mathbb{F}} &\leq \sqrt{\frac{3}{\ln k}} \epsilon e \cdot \left\| \mathbf{\Sigma}_2^{1/2} \right\|_{\mathbb{F}} + \frac{2e^2}{\sqrt{k}} \epsilon \cdot \left\| \mathbf{\Sigma}_2^{1/2} \right\|_2 \\ &= \epsilon e \sqrt{\frac{3}{\ln k}} \left\| \mathbf{\Sigma}_2 \right\|_{\star} + \frac{2e^2}{\sqrt{k}} \epsilon \cdot \left\| \mathbf{\Sigma}_2 \right\|_2^{1/2} \end{aligned}$$

with the same probability. Likewise,

$$\begin{aligned} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\mathbb{F}}^2 &\leq \left( \epsilon e \sqrt{\frac{3}{\ln k}} \left\| \mathbf{\Sigma}_2 \right\|_{\star} + \frac{2e^2}{\sqrt{k}} \epsilon \cdot \left\| \mathbf{\Sigma}_2 \right\|_2^{1/2} \right)^2 \\ &\leq \frac{6}{\ln k} \epsilon^2 e^2 \cdot \left\| \mathbf{\Sigma}_2 \right\|_{\star} + \frac{8e^4}{k} \epsilon^2 \cdot \left\| \mathbf{\Sigma}_2 \right\|_2 \end{aligned}$$

with the same probability.

These estimates used in Theorems 2 and 4 yield the stated spectral and trace norm bounds. To obtain the corresponding Frobenius norm bound, define the quantities

$$\begin{aligned} G_1 &= \left( \frac{12e^2}{\ln k} + \frac{16e^4}{k} \right) \epsilon^2 & G_3 &= 3e^2 \frac{\epsilon^2}{\ln k} \\ G_2 &= 4e^4 \frac{\epsilon^2}{k \ln k} & G_4 &= 4e^4 \frac{\epsilon^2}{k} \end{aligned}$$



By Theorem 3 and our estimates for  $\left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_2$  and  $\left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_F$ ,

$$\begin{aligned}
 \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_F &\leq \left\| \Sigma_2 \right\|_F + \gamma^{q-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_2 \cdot \left( \sqrt{2 \operatorname{Tr}(\Sigma_2)} + \gamma^{q-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_F \right) \\
 &\leq \left\| \Sigma_2 \right\|_F + \gamma^{q-1} (G_1 \left\| \Sigma_2 \right\|_2 + G_2 \left\| \Sigma_2 \right\|_\star)^{1/2} \times \\
 &\quad \left( \sqrt{2 \operatorname{Tr}(\Sigma_2)} + \gamma^{q-1} \sqrt{G_3 \left\| \Sigma_2 \right\|_\star} + \gamma^{q-1} \sqrt{G_4 \left\| \Sigma_2 \right\|_2} \right) \\
 &\leq \left\| \Sigma_2 \right\|_F + \left( \gamma^{q-1} \sqrt{2G_1} + \gamma^{2q-2} (\sqrt{G_1 G_3} + \sqrt{G_2 G_4}) \right) \cdot \sqrt{\left\| \Sigma_2 \right\|_2 \left\| \Sigma_2 \right\|_\star} \\
 &\quad + \left( \gamma^{q-1} \sqrt{2G_2} + \gamma^{2q-2} \sqrt{G_2 G_3} \right) \cdot \left\| \Sigma_2 \right\|_\star \\
 &\quad + \gamma^{2q-2} \sqrt{G_1 G_4} \left\| \Sigma_2 \right\|_2.
 \end{aligned} \tag{21}$$

The following estimates hold for the coefficients in this inequality:

$$\begin{aligned}
 \sqrt{2G_1} &\leq \left( \frac{42}{\sqrt{k}} + \frac{14}{\sqrt{\ln k}} \right) \epsilon & \sqrt{G_1 G_3} &\leq \left( \frac{45}{\ln k} + \frac{140}{\sqrt{k \ln k}} \right) \epsilon^2 \\
 \sqrt{G_2 G_4} &\leq \frac{219}{k \sqrt{\ln k}} \epsilon^2 & \sqrt{2G_2} &\leq 21 \frac{\epsilon}{\sqrt{k \ln k}} \\
 \sqrt{G_2 G_3} &\leq 70 \frac{\epsilon^2}{\sqrt{k \ln k}} & \sqrt{G_1 G_4} &\leq \left( \frac{140}{\sqrt{k \ln k}} + \frac{437}{k} \right) \epsilon^2.
 \end{aligned}$$

The Frobenius norm bound follows from using these estimates in Equation (21) and grouping terms appropriately:

$$\begin{aligned}
 \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_F &\leq \left\| \Sigma_2 \right\|_F + \left[ \gamma^{q-1} \epsilon \left( \frac{42}{\sqrt{k}} + \frac{14}{\sqrt{\ln k}} \right) \right. \\
 &\quad \left. + \gamma^{2q-2} \epsilon^2 \left( \frac{45}{\ln k} + \frac{140}{\sqrt{k \ln k}} + \frac{219}{k \sqrt{\ln k}} \right) \right] \sqrt{\left\| \Sigma_2 \right\|_2 \left\| \Sigma_2 \right\|_\star} \\
 &\quad + \left( 21 \gamma^{q-1} \frac{\epsilon}{\sqrt{k \ln k}} + 70 \gamma^{2q-2} \frac{\epsilon^2}{\sqrt{k \ln k}} \right) \cdot \left\| \Sigma_2 \right\|_\star \\
 &\quad + \gamma^{2q-2} \epsilon^2 \left( \frac{140}{\sqrt{k \ln k}} + \frac{437}{k} \right) \left\| \Sigma_2 \right\|_2.
 \end{aligned}$$

■

#### 4.2.4 SAMPLING COLUMNS UNIFORMLY AT RANDOM

Here, the columns of  $\mathbf{A}$  are sampled uniformly at random (with or without replacement). Such uniformly-at-random column sampling only makes sense when the leverage scores of the top  $k$ -dimensional invariant subspace of the matrix are sufficiently uniform that no column is significantly more informative than the others. For this case, we can prove the following.

**Lemma 8** *Let  $\mathbf{A}$  be an  $n \times n$  SPSD matrix,  $q$  be a positive integer, and  $\mathbf{S}$  be a sampling matrix of size  $n \times \ell$  corresponding to sampling the columns of  $\mathbf{A}$  uniformly at random (with*

or without replacement). Let  $\mu$  denote the coherence of the top  $k$ -dimensional eigenspace of  $\mathbf{A}$  and fix a failure probability  $\delta \in (0, 1)$  and accuracy factor  $\epsilon \in (0, 1)$ . Define

$$\gamma = \frac{\lambda_{k+1}(\mathbf{A})}{\lambda_k(\mathbf{A})}.$$

If  $\ell \geq 2\mu\epsilon^{-2}k \ln(k/\delta)$ , then, when  $\mathbf{C} = \mathbf{A}^q \mathbf{S}$  and  $\mathbf{W} = \mathbf{S}^T \mathbf{A}^{2q-1} \mathbf{S}$ , the corresponding low-rank SPSP approximation satisfies

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_2 &\leq \left( 1 + \left( \frac{n}{(1-\epsilon)\ell} \right)^{1/(2q-1)} \right) \|\mathbf{A} - \mathbf{A}_k\|_2, \\ \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\text{F}} &\leq \|\mathbf{A} - \mathbf{A}_k\|_{\text{F}} + \left( \gamma^{q-1} \frac{\sqrt{2}}{\delta \sqrt{1-\epsilon}} + \frac{\gamma^{2q-2}}{(1-\epsilon)\delta^2} \right) \|\mathbf{A} - \mathbf{A}_k\|_{\star}, \text{ and} \\ \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\star} &\leq \left( 1 + \frac{\gamma^{2q-2}}{\delta^2(1-\epsilon)} \right) \|\mathbf{A} - \mathbf{A}_k\|_{\star}, \end{aligned}$$

simultaneously with probability at least  $1 - 3\delta$ .

**Proof** In (Gittens, 2012), it is shown that

$$\left\| \mathbf{\Omega}_1^\dagger \right\|_2^2 \leq \frac{n}{(1-\epsilon)\ell}$$

with probability at least  $1 - \delta$  when  $\ell$  satisfies the stated bound. Observe that  $\|\mathbf{\Omega}_2\|_2 \leq \|\mathbf{U}_2\|_2 \|\mathbf{S}\|_2 \leq 1$ , so that

$$\left\| \mathbf{\Sigma}_2^{q-1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_2^2 \leq \left\| \mathbf{\Sigma}_2^{q-1/2} \right\|_2^2 \left\| \mathbf{\Omega}_1^\dagger \right\|_2^2 \leq \|\mathbf{\Sigma}_2\|_2^{2q-1} \frac{n}{(1-\epsilon)\ell}$$

with probability at least  $1 - \delta$ . Also,

$$\left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \mathbf{\Omega}_1^\dagger \right\|_{\text{F}} \leq \sqrt{\frac{n}{(1-\epsilon)\ell}} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \right\|_{\text{F}} \quad (22)$$

with at least the same probability. Observe that since  $\mathbf{S}$  selects  $\ell$  columns uniformly at random,

$$\mathbb{E} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \right\|_{\text{F}}^2 = \mathbb{E} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{U}_2^T \mathbf{S} \right\|_{\text{F}}^2 = \sum_{i=1}^{\ell} \mathbb{E} \|\mathbf{x}_i\|^2,$$

where the summands  $\mathbf{x}_i$  are distributed uniformly at random over the columns of  $\mathbf{\Sigma}_2^{1/2} \mathbf{U}_2^T$ . Regardless of whether  $\mathbf{S}$  selects the columns with replacement or without replacement, the summands all have the same expectation:

$$\mathbb{E} \|\mathbf{x}_i\|^2 = \frac{1}{n} \sum_{j=1}^n \left\| (\mathbf{\Sigma}_2^{1/2} \mathbf{U}_2^T)^j \right\|^2 = \frac{1}{n} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{U}_2^T \right\|_{\text{F}}^2 = \frac{1}{n} \left\| \mathbf{\Sigma}_2^{1/2} \right\|_{\text{F}}^2 = \frac{1}{n} \|\mathbf{\Sigma}_2\|_{\star}.$$

Consequently,

$$\mathbb{E} \left\| \mathbf{\Sigma}_2^{1/2} \mathbf{\Omega}_2 \right\|_{\text{F}}^2 = \frac{\ell}{n} \|\mathbf{\Sigma}_2\|_{\star},$$

so by Jensen's inequality

$$\mathbb{E} \left\| \Sigma_2^{1/2} \Omega_2 \right\|_{\text{F}} \leq \left( \mathbb{E} \left\| \Sigma_2^{1/2} \Omega_2 \right\|_{\text{F}}^2 \right)^{1/2} = \sqrt{\frac{\ell}{n}} \|\Sigma_2\|_{\star}.$$

Now applying Markov's inequality to (22), we see that

$$\left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}} \leq \frac{1}{\delta} \sqrt{\frac{1}{(1-\epsilon)}} \|\Sigma_2\|_{\star}$$

with probability at least  $1 - 2\delta$ . Thus, we also know that

$$\left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}}^2 \leq \frac{1}{(1-\epsilon)\delta^2} \|\Sigma_2\|_{\star},$$

also with probability at least  $1 - 2\delta$ . These estimates used in Theorems 2 and 4 yield the stated spectral and trace norm bounds.

To obtain the Frobenius norm bound, observe that Theorem 3 implies

$$\begin{aligned} \left\| \mathbf{A} - \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^T \right\|_{\text{F}} &\leq \|\Sigma_2\|_{\text{F}} + \gamma^{p-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}} \left( \sqrt{2 \text{Tr}(\Sigma_2)} + \gamma^{p-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}} \right) \\ &\leq \|\Sigma_2\|_{\text{F}} + \gamma^{p-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}} \left( \sqrt{2 \text{Tr}(\Sigma_2)} + \gamma^{p-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}} \right) \\ &\leq \|\Sigma_2\|_{\text{F}} + \gamma^{2p-2} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}}^2 + \gamma^{p-1} \left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}} \sqrt{2 \text{Tr}(\Sigma_2)}. \end{aligned}$$

Now substitute our estimate for  $\left\| \Sigma_2^{1/2} \Omega_2 \Omega_1^\dagger \right\|_{\text{F}}^2$  to obtain the stated Frobenius norm bound. ■

*Remark.* As with previous bounds for uniform sampling, (e.g., Kumar et al., 2012; Gittens, 2012), these results for uniform sampling are much weaker than our bounds from the previous subsections, since the sampling complexity depends on the coherence of the input matrix. When the matrix has small coherence, however, these bounds are similar to the bounds derived from the leverage-based sampling probabilities. Recall that, by the algorithm of Drineas et al. (2012), the coherence of an arbitrary input matrix can be computed in roughly the time it takes to perform a random projection on the input matrix.

## 5. Discussion and Conclusion

We have presented a unified approach to a large class of low-rank approximations of Laplacian and kernel matrices that arise in machine learning and data analysis applications. In doing so, we have provided qualitatively-improved worst-case theory and clarified the performance of these algorithms in practical settings. Our theoretical and empirical results suggest several obvious directions for future work.

In general, our empirical evaluation demonstrates that obtaining moderately high-quality low-rank approximations, as measured by minimizing the reconstruction error, depends in complicated ways on the spectral decay, the leverage score structure, the eigenvalue gaps in

relevant parts of the spectrum, etc. (Ironically, our empirical evaluation also demonstrates that *all* the sketches considered are reasonably-effective at approximating both sparse and dense, and both low-rank and high-rank matrices which arise in practice. That is, with only roughly  $O(k)$  measurements, the spectral, Frobenius, and trace approximation errors stay within a small multiplicative factor of around 3 of the optimal rank- $k$  approximation errors. The reason for this is that matrices for which uniform sampling is least appropriate tend to be those which are least well-approximated by low-rank matrices, meaning that the residual error is much larger.) Thus, *e.g.*, depending on whether one is interested in  $\ell$  being slightly larger or much larger than  $k$ , leverage-based sampling or a random projection might be most appropriate; and, more generally, an ensemble-based method that draws complementary strengths from each of these methods might be best.

In addition, we should note that, in situations where one is concerned with the quality of approximation of the actual eigenspaces, one desires both a small spectral norm error (because by the Davis–Kahan  $\sin \Theta$  theorem and similar perturbation results, this would imply that the range space of the sketch effectively captures the top  $k$ -dimensional eigenspace of  $\mathbf{A}$ ) as well as to use as few samples as possible (because one prefers to approximate the top  $k$ -dimension eigenspace of  $\mathbf{A}$  with as close to a  $k$ -dimensional subspace as possible). Our results suggest that the leverage score probabilities supply the best sampling scheme for balancing these two competing objectives.

More generally, although our empirical evaluation consists of *random* sampling and *random* projection algorithms, our theoretical analysis clearly decouples the randomness in the algorithm from the structural heterogeneities in the Euclidean vector space that are responsible for the poor performance of uniform sampling algorithms. Thus, if those structural conditions can be satisfied with a deterministic algorithm, an iterative algorithm, or any other method, then one can certify (after running the algorithm) that good approximation guarantees hold for particular input matrices in less time than is required for general matrices. Moreover, this structural decomposition suggests greedy heuristics—*e.g.*, greedily keep some number of columns according to approximate statistical leverage scores and “residualize.” In our experience, a procedure of this form often performs quite well in practice, although theoretical guarantees tend to be much weaker; and thus we expect that, when coupled with our results, such procedures will perform quite well in practice in many medium-scale and large-scale machine learning applications.

## Acknowledgments

AG would like to acknowledge the support, under the auspice of Joel Tropp, of ONR awards N00014-08-1-0883 and N00014-11-1-0025, AFOSR award FA9550-09-1-0643, and a Sloan Fellowship; and MM would like to acknowledge a grant from the Defense Advanced Research Projects Agency.

## References

N. Arcolano and P. J. Wolfe. Nyström Approximation of Wishart Matrices. In *Proceedings of the 2010 IEEE International Conference on Acoustics Speech and Signal Processing*,

- pages 3606–3609, 2010.
- A. Asuncion and D. J. Newman. UCI Machine Learning Repository, November 2012. URL <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- H. Avron, P. Maymounkov, and S. Toledo. Blendenpik: Supercharging LAPACK’s Least-squares Solver. *SIAM Journal on Scientific Computing*, 32:1217–1236, 2010.
- F. Bach. Sharp Analysis of Low-rank Kernel Matrix Approximations. In *Proceedings of the International Conference on Learning Theory (COLT)*, 2013.
- F.R. Bach and M.I. Jordan. Predictive Low-rank Decomposition for Kernel Methods. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 33–40, 2005.
- A. Banerjee, D. Dunson, and S. Tokdar. Efficient Gaussian Process Regression for Large Data Sets. *Biometrika*, 100:75–89, 2012.
- M.-A. Belabbas and P. J. Wolfe. Fast Low-Rank Approximation for Covariance Matrices. In *Second IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing*, pages 293–296, 2007a.
- M.-A. Belabbas and P. J. Wolfe. On Sparse Representations of Linear Operators and the Approximation of Matrix Products. In *Proceedings of the 42nd Annual Conference on Information Sciences and Systems*, pages 258–263, 2008.
- M.-A. Belabbas and P. J. Wolfe. On Landmark Selection and Sampling in High-dimensional Data Analysis. *Philosophical Transactions of the Royal Society, Series A*, 367:4295–4312, 2009a.
- M.-A. Belabbas and P. J. Wolfe. Spectral Methods in Machine Learning and New Strategies for Very Large Datasets. *Proc. Natl. Acad. Sci. USA*, 106:369–374, 2009b.
- M.-A. Belabbas and P.J. Wolfe. On the Approximation of Matrix Products and Positive Definite Matrices. Technical report, 2007b. Preprint: arXiv:0707.4448 (2007).
- E. Bingham and H. Mannila. Random Projection in Dimensionality Reduction: Applications to Image and Text Data. In *Proceedings of the 7th Annual ACM SIGKDD Conference*, pages 245–250, 2001.
- C. Boutsidis and A. Gittens. Improved Matrix Algorithms Via the Subsampled Randomized Hadamard Transform. *SIAM Journal of Matrix Analysis and Its Applications*, 34:1301–1340, 2013.
- C. Boutsidis, M.W. Mahoney, and P. Drineas. An Improved Approximation Algorithm for the Column Subset Selection Problem. In *Proceedings of the 20th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 968–977, 2009.
- J. Chiu and L. Demanet. Sublinear Randomized Algorithms for Skeleton Decompositions. *SIAM Journal on Matrix Analysis and Its Applications*, 34:1361–1383, 2013.

- P. I. Corke. A Robotics Toolbox for MATLAB. *IEEE Robotics and Automation Magazine*, 3:24–32, 1996.
- C. Cortes, M. Mohri, and A. Talwalkar. On the Impact of Kernel Approximation on Learning Accuracy. In *Proceedings of the 13th International Workshop on Artificial Intelligence and Statistics*, 2010.
- P. Drineas and M.W. Mahoney. On the Nyström Method for Approximating a Gram Matrix for Improved Kernel-Based Learning. *Journal of Machine Learning Research*, 6: 2153–2175, 2005.
- P. Drineas, R. Kannan, and M.W. Mahoney. Fast Monte Carlo Algorithms for Matrices I: Approximating Matrix Multiplication. *SIAM Journal on Computing*, 36:132–157, 2006.
- P. Drineas, M.W. Mahoney, and S. Muthukrishnan. Relative-error CUR Matrix Decompositions. *SIAM Journal on Matrix Analysis and Applications*, 30:844–881, 2008.
- P. Drineas, M.W. Mahoney, S. Muthukrishnan, and T. Sarlós. Faster Least Squares Approximation. *Numerische Mathematik*, 117(2):219–249, 2010.
- P. Drineas, M. Magdon-Ismael, M. W. Mahoney, and D. P. Woodruff. Fast Approximation of Matrix Coherence and Statistical Leverage. *Journal of Machine Learning Research*, 13:3475–3506, 2012.
- A. K. Farahat, A. Ghodsi, and M. S. Kamel. A Novel Greedy Algorithm for Nyström Approximation. In *Proceedings of the 14th International Workshop on Artificial Intelligence and Statistics*, 2011.
- C. Fowlkes, S. Belongie, F. Chung, and J. Malik. Spectral Grouping Using the Nyström Method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(2):214–225, 2004.
- D. Fradkin and D. Madigan. Experiments with Random Projections for Machine Learning. In *Proceedings of the 9th Annual ACM SIGKDD Conference*, pages 517–522, 2003.
- M. Genton. Classes of Kernels for Machine Learning: a Statistics Perspective. *J. Mach. Learn. Res.*, 2:299–312, 2002.
- A. Gittens. The Spectral Norm Error of the Naïve Nyström Extension. Technical report, California Institute of Technology, 2012. Preprint: arXiv:1110.5305 (2011).
- A. Gittens and M. W. Mahoney. Revisiting the Nyström Method for Improved Large-scale Machine Learning. Technical report, 2013. Tech Report: arXiv:1303.1849.
- G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, 3rd edition, 1996.
- A. M. Gustafson, E. S. Snitkin, S. C. J. Parker, C. DeLisi, and S. Kasif. Towards the Identification of Essential Genes Using Targeted Genome Sequencing and Comparative Analysis. *BMC Genomics*, 7:265, 2006.

- I. Guyon, S. R. Gunn, A. Ben-Hur, and G. Dror. Result Analysis of the NIPS 2003 Feature Selection Challenge. In *Advances in Neural Information Processing Systems 17*. MIT Press, 2005.
- N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions. *SIAM Review*, 53(2):217–288, 2011.
- D. Homrighausen and D. J. McDonald. Spectral Approximations in Machine Learning. Technical report, 2011. Preprint: arXiv:1107.4340 (2011).
- R. Jin, T. Yang, M. Mahdavi, Y.-F. Li, and Z.-H. Zhou. Improved Bound for the Nyström’s Method and its Application to Kernel Classification. *IEEE Information Theory*, 59:6939–6949, 2013. Preprint: arXiv:1111.2262 (2011).
- B. Klimt and Y. Yang. The Enron Corpus: a New Dataset for Email Classification Research. In *Proceedings of the 15th European Conference on Machine Learning*, pages 217–226, 2004.
- S. Kumar, M. Mohri, and A. Talwalkar. On Sampling-based Approximate Spectral Decomposition. In *Proceedings of the 26th International Conference on Machine Learning*, pages 553–560, 2009a.
- S. Kumar, M. Mohri, and A. Talwalkar. Ensemble Nyström Method. In *Annual Advances in Neural Information Processing Systems 22: Proceedings of the 2009 Conference*, 2009b.
- S. Kumar, M. Mohri, and A. Talwalkar. Sampling Techniques for the Nyström Method. In *Proceedings of the 12th Tenth International Workshop on Artificial Intelligence and Statistics*, pages 304–311, 2009c.
- S. Kumar, M. Mohri, and A. Talwalkar. Sampling Methods for the Nyström Method. *Journal of Machine Learning Research*, 13:981–1006, 2012.
- J. Leskovec, J. Kleinberg, and C. Faloutsos. Graph Evolution: Densification and Shrinking Diameters. *ACM Transactions on Knowledge Discovery from Data*, 1, 2007.
- M. Li, J.T. Kwok, and B.-L. Lu. Making Large-Scale Nyström Approximation Possible. In *Proceedings of the 27th International Conference on Machine Learning*, pages 631–638, 2010.
- S. Liu, J. Zhang, and K. Sun. Learning Low-rank Kernel Matrices with Column-based Methods. *Communications in Statistics—Simulation and Computation*, 39(7):1485–1498, 2010.
- P. Ma, M. W. Mahoney, and B. Yu. A Statistical Perspective on Algorithmic Leveraging. In *Proceedings of the 31th International Conference on Machine Learning*, 2014.
- P. Machart, T. Peel, S. Anthoine, L. Ralaivola, and H. Glotin. Stochastic Low-Rank Kernel Learning for Regression. In *Proceedings of the 28th International Conference on Machine Learning*, pages 969–976, 2011.

- L. Mackey, A. Talwalkar, and M. I. Jordan. Divide-and-conquer Matrix Factorization. In *Annual Advances in Neural Information Processing Systems 24: Proceedings of the 2011 Conference*, 2011a.
- L. Mackey, A. Talwalkar, and M. I. Jordan. Divide-and-conquer Matrix Factorization. In *NIPS*, 2011b. Preprint: arXiv:1107.0789 (2011).
- M. W. Mahoney. *Randomized Algorithms for Matrices and Data*. Foundations and Trends in Machine Learning. NOW Publishers, Boston, 2011. Also available at: arXiv:1104.5557.
- M. W. Mahoney. Algorithmic and Statistical Perspectives on Large-Scale Data Analysis. In U. Naumann and O. Schenk, editors, *Combinatorial Scientific Computing*, Chapman & Hall/CRC Computational Science. CRC Press, 2012.
- M.W. Mahoney and P. Drineas. CUR Matrix Decompositions for Improved Data Analysis. *Proc. Natl. Acad. Sci. USA*, 106:697–702, 2009.
- P.-G. Martinsson, V. Rokhlin, and M. Tygert. A Randomized Algorithm for the Decomposition of Matrices. *Applied and Computational Harmonic Analysis*, 30:47–68, 2011.
- X. Meng, M. A. Saunders, and M. W. Mahoney. LSRN: a Parallel Iterative Solver for Strongly Over- or Under-Determined Systems. *SIAM Journal on Scientific Computing*, 36(2):C95–C118, 2014.
- M. Mohri and A. Talwalkar. Can Matrix Coherence be Efficiently and Accurately Estimated? In *Proceedings of the 14th International Workshop on Artificial Intelligence and Statistics*, 2011.
- T. O. Nielsen, R. B. West, S. C. Linn, O. Alter, M. A. Knowling, J. X. O’Connell, S. Zhu, M. Fero, G. Sherlock, J. R. Pollack, P. O. Brown, D. Botstein, and M. van de Rijn. Molecular Characterisation of Soft Tissue Tumours: a Gene Expression Study. *The Lancet*, 359:1301–1307, 2002.
- P. Parker, P. J. Wolfe, and V. Tarok. A Signal Processing Application of Randomized Low-rank Approximations. In *Proceedings of the 13th IEEE Workshop on Statistical Signal Processing*, pages 345–350, 2005.
- P. Paschou, E. Ziv, E.G. Burchard, S. Choudhry, W. Rodriguez-Cintron, M.W. Mahoney, and P. Drineas. PCA-Correlated SNPs for Structure Identification in Worldwide Human Populations. *PLoS Genetics*, 3:1672–1686, 2007.
- V. Rokhlin, A. Szlam, and M. Tygert. A Randomized Algorithm for Principal Component Analysis. *SIAM Journal on Matrix Analysis and Applications*, 31(3):1100–1124, 2009.
- B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001.
- D. N. Spendley and P. J. Wolfe. Adaptive Beamforming Using Fast Low-rank Covariance Matrix Approximations. In *Proceedings of the IEEE Radar Conference*, pages 1–5, 2008.



- A. Talwalkar and A. Rostamizadeh. Matrix Coherence and the Nyström Method. In *Proceedings of the 26th Conference in Uncertainty in Artificial Intelligence*, 2010.
- A. Talwalkar, S. Kumar, and H. Rowley. Large-scale Manifold Learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008.
- J. A. Tropp. Improved Analysis of the Subsampled Randomized Hadamard Transform. *Adv. Adapt. Data Anal.*, 3(1-2):115–126, 2011.
- S. Venkatasubramanian and Q. Wang. The Johnson–Lindenstrauss Transform: An Empirical Study. In *ALENEX11: Workshop on Algorithms Engineering and Experimentation*, pages 164–173, 2011.
- S. Wang and Z. Zhang. Improving CUR Matrix Decomposition and Nyström Approximation via Adaptive Sampling. *Journal of Machine Learning Research*, 14:2549–2589, 2013. Preprint: arXiv:1303.4207 (2013).
- C.K.I. Williams and M. Seeger. Using the Nyström Method to Speed Up Kernel Machines. In *Annual Advances in Neural Information Processing Systems 13: Proceedings of the 2000 Conference*, pages 682–688, 2001.
- C.K.I. Williams, C.E. Rasmussen, A. Schwaighofer, and V. Tresp. Observations on the Nyström Method for Gaussian Process Prediction. Technical report, University of Edinburgh, 2002.
- F. Woolfe, E. Liberty, V. Rokhlin, and M. Tygert. A Fast Randomized Algorithm for the Approximation of Matrices. *Applied and Computational Harmonic Analysis*, 25(3):335–366, 2008.
- C.-W. Yip, M. W. Mahoney, A. S. Szalay, I. Csabai, T. Budavári, R. F. G. Wyse, and L. Dobos. Objective Identification of Informative Wavelength Regions in Galaxy Spectra. *The Astronomical Journal*, 147(5):110, 2014.
- K. Zhang and J. T. Kwok. Density-weighted Nyström Method for Computing Large Kernel Eigensystems. *Neural Computation*, 21(1):121–146, 2009.
- K. Zhang and J. T. Kwok. Clustered Nyström Method for Large Scale Manifold Learning and Dimension Reduction. *IEEE Transactions on Neural Networks*, 21(10):1576–1587, 2010.
- K. Zhang, I.W. Tsang, and J.T. Kwok. Improved Nyström Low-rank Approximation and Error Analysis. In *Proceedings of the 25th International Conference on Machine Learning*, pages 1232–1239, 2008.