# Quantum-Inspired Algorithms from Randomized Numerical Linear Algebra

Nadiia Chepurko <sup>1</sup> Kenneth L. Clarkson <sup>2</sup> Lior Horesh <sup>2</sup> Honghao Lin <sup>3</sup> David P. Woodruff <sup>3</sup>

#### **Abstract**

We create classical (non-quantum) dynamic data structures supporting queries for recommender systems and least-squares regression that are comparable to their quantum analogues. Dequantizing such algorithms has received a flurry of attention in recent years; we obtain sharper bounds for these problems. More significantly, we achieve these improvements by arguing that the previous quantum-inspired algorithms for these problems are doing leverage or ridge-leverage score sampling in disguise; these are powerful and standard techniques in randomized numerical linear algebra. With this recognition, we are able to employ the large body of work in numerical linear algebra to obtain algorithms for these problems that are simpler or faster (or both) than existing approaches. Our experiments demonstrate that the proposed data structures also work well on real-world datasets.

#### 1. Introduction

In recent years, quantum algorithms for various problems in numerical linear algebra have been proposed, with applications including least-squares regression and recommender systems (Harrow et al., 2009; Lloyd et al., 2016; Rebentrost et al., 2014; Gilyén et al., 2019; Zhao et al., 2019; Brandão et al., 2019; van Apeldoorn & Gilyén, 2019; Lloyd et al., 2014; Cong & Duan, 2016; Berry et al., 2015). Some of these algorithms have the striking property that their running times do not depend on the input size. That is, for a given matrix  $A \in \mathbb{R}^{n \times d}$  with  $\mathtt{nnz}(A)$  nonzero entries, the

Proceedings of the 39<sup>th</sup> International Conference on Machine Learning, Baltimore, Maryland, USA, PMLR 162, 2022. Copyright 2022 by the author(s).

running times for these proposed quantum algorithms are at most polylogarithmic in n and d, and polynomial in other parameters of A, such as  $\operatorname{rank}(A)$ , the condition number  $\kappa(A)$ , or Frobenius norm  $\|A\|_F$ .

However, as observed by Tang (Tang, 2019) and others, there is a catch: these quantum algorithms depend on a particular input representation of A, which is a simple data structure that allows A to be employed for quick preparation of a quantum state suitable for further quantum computations. This data structure, which is a collection of weighted complete binary trees, also supports rapid weighted random sampling of A, for example, sampling the rows of A with probability proportional to their squared Euclidean lengths. So, if an "apples to apples" comparison of quantum to classical computation is to be made, it is reasonable to ask what can be accomplished in the classical realm using the sampling that the given data structure supports.

A recent line of work analyzes the speedups of these quantum algorithms by developing classical counterparts that exploit these restrictive input and output assumptions, and shows that previous quantum algorithms do not give an exponential speedup. In this setting, it has recently been shown that sublinear time is sufficient for least-squares regression using a low-rank design matrix A (Gilyén et al., 2018; Chia et al., 2018), for computing a low-rank approximation to input matrix A (Tang, 2019), and for solving ridge regression problems (Gilyén et al., 2020), using classical (non-quantum) methods, assuming the data structure of trees has already been constructed. Further, the results obtained in (Tang, 2019; Gilyén et al., 2018; 2020) serve as appropriate comparisons of the power of quantum to classical computing, due to their novel input-output model: data structures are input, then sublinear-time computations are done, yielding data structures as output.

The simple weighted-sampling data structure used in these works to represent the input can be efficiently constructed and stored: it uses  $O(\mathtt{nnz}(A))$  space, with a small constant overhead, and requires  $O(\mathtt{nnz}(A))$  time to construct, in the static case where the matrix A is given in its entirety, and can support updates and queries to individual entries of A in  $O(\log(nd))$  time. However, the existing reported sublinear bounds are high-degree polynomials in

<sup>&</sup>lt;sup>1</sup>Department of Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA <sup>2</sup>IBM Research, USA <sup>3</sup>Department of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA. Correspondence to: Kenneth Clarkson <klclarks@us.ibm.com>, Honghao Lin <honghaol@andrew.cmu.edu>, David Woodruff <dwoodruf@andrew.cmu.edu>.

the parameters involved: for instance, the sublinear term in the running time for low-rank least-squares regression is  $\tilde{O}({\tt rank}(A)^6\|A\|_F^6\kappa(A)^{16}/\varepsilon^6)$ ; see also more recent work for ridge regression (Gilyén et al., 2020). This combination of features raises the following question:

Question 1: Can the sublinear terms in the running time be reduced significantly while preserving the leading order dependence of O(nnz(A)) and  $O(\log(nd))$  per update (dynamic)?

Perhaps a question of greater importance is the connection between quantum-inspired algorithms and the vast body of work in randomized numerical linear algebra: see the surveys (Kannan & Vempala, 2009; Mahoney, 2011; Woodruff, 2014). There are a large number of randomized algorithms based on sampling and sketching techniques for problems in linear algebra. Yet prior to our work, none of the quantum-inspired algorithms, which are sampling-based, have discussed the connection to leverage scores, for example, which are a powerful and standard tool.

Question 2: Can the large body of work in randomized numerical linear algebra be applied effectively in the setting of quantum-inspired algorithms?

#### 1.1. Our Results

We answer both of the questions above affirmatively. In fact, we answer Question 1 by answering Question 2. Namely, we obtain significant improvements in the sublinear terms, and our analysis relies on simulating *leverage score sampling* and *ridge leverage score sampling*, using the aforementioned data structure to sample rows proportional to squared Euclidean norms. Additionally, we empirically demonstrate the speedup we achieve on real-world and synthetic datasets (see Section 6).

Connection to Classical Linear Algebra and Dynamic Data Structures. The work on quantum-inspired algorithms builds data structures for sampling according to the squared row and column lengths of a matrix. This is also a common technique in randomized numerical linear algebra - see the recent survey on length-squared sampling by Kannan and Vempala (Kannan & Vempala, 2017). However, it is well-known that leverage score sampling often gives stronger guarantees than length-squared sampling; leverage score sampling was pioneered in the algorithms community in (Drineas et al., 2006), and made efficient in (Drineas et al., 2012) (see also analogous prior work in the  $\ell_1$  setting, starting with (Clarkson, 2005)).

Given an  $n \times d$  matrix A, with  $n \ge d$ , its (row) leverage scores are the squared row norms of U, where U is an orthonormal basis with the same column span as A. One can show that any choice of basis gives the same scores. Writing

 $A = U\Sigma V^T$  in its thin singular value decomposition (SVD), and letting  $A_{i,*}$  and  $U_{i,*}$  denote the i-th rows of A and U respectively, we see that  $\|A_{i,*}\| = \|U_{i,*}\Sigma\|$ . Consequently, letting  $k = \operatorname{rank}(A)$ , and with  $\sigma_1$  and  $\sigma_k$  denoting the maximum and minimum non-zero singular values of A, we have  $\|A_{i,*}\| \geq \|U_{i,*}\|\sigma_k(A)$ , and  $\|A_{i,*}\| \leq \|U_{i,*}\|\sigma_1(A)$ .

Thus, sampling according to the squared row norms of A is equivalent to sampling from a distribution with ratio distance at most  $\kappa^2(A) = \frac{\sigma_1(A)^2}{\sigma_k(A)^2}$  from the leverage score distribution, that is, sampling a row with probability proportional to its leverage score. This is crucial, as it implies using standard arguments (see, e.g., (Woodruff, 2014) for a survey) that if we oversample by a factor of  $\kappa^2(A)$ , then we obtain the same guarantees for various problems that leverage score sampling achieves. Notice that the running times of quantum-inspired algorithms, e.g., the aforementioned  $\tilde{O}(\operatorname{rank}(A)^6 \|A\|_F^6 \kappa(A)^{16}/\varepsilon^6)$ time for regression of (Gilyén et al., 2018) and the  $\tilde{O}(\|A\|_F^8 \kappa(A)^2/(\sigma_{min}(A)^6 \varepsilon^4))$  time for regression of (Gilyén et al., 2020), both take a number of squared-length samples of A depending on  $\kappa(A)$ , and thus are implicitly doing leverage score sampling, or in the case of ridge regression, ridge leverage score sampling.

Given the connection above, we focus on two central problems in machine learning and numerical linear algebra, ridge regression (Problem 1.1) and low rank approximation (Problem 1.2). We show how to obtain simpler algorithms and analysis than those in the quantum-inspired literature by using existing approximate matrix product and subspace embedding guarantees of leverage score sampling. In addition to improved bounds, our analysis de-mystifies what the rather involved  $\ell_2$ -sampling arguments of quantum-inspired work are doing, and decreases the gap between quantum and classical algorithms for machine learning problems. We begin by formally defining ridge regression and low-rank approximation, and the dynamic data structure model we focus on

Problem 1.1 (Ridge Regression). Given an  $n \times d$  matrix A,  $n \times d'$  matrix B and a ridge parameter  $\lambda \geq 0$ , the ridge regression problem is defined as follows:

$$\min_{X \in \mathbb{R}^{d \times d'}} \|AX - B\|_F^2 + \lambda \|X\|_F^2,$$

where  $\|\cdot\|_F^2$  denotes the sum-of-squares of entries.

Problem 1.2 (Low-Rank Approximation). Given an  $n \times d$  matrix A and a rank parameter  $k \in [d]$ , the low-rank approximation problem is defined as follows:

$$\min_{X \in \mathbb{R}^{n \times d}: \mathrm{rank}(X) = k} \|A - X\|_F^2.$$

**Definition 1.3** (Dynamic Data Structure Model). Given an  $n \times d$  matrix A, the dynamic data structure supports the

following operations in  $O(\log(nd))$  time: (a) sample row  $A_{i,*}$  with probability  $\|A_{i,*}\|_2^2/\|A\|_F^2$ , (b) sample entry j in row i with probability  $A_{i,j}^2/\|A_{i,*}\|_2^2$  and (c) output the (i,j)-th entry of A.

We note that in this input model, reading the entire matrix would be prohibitive and the algorithms can only access the matrix through the weighted sampling data structure.

We now describe our concrete results in more detail. At a high level, our algorithm for ridge regression, Algorithm 2, does the following: sample a subset of the rows of A via length-squared sampling, and take a length-squared sample of the columns of that subset. Then, solve a linear system on the resulting small submatrix using the conjugate gradient method. Our analysis of this algorithm results in the following theorem. (Some of the (standard) matrix notation used is given in Section 3.)

The residual error is bounded in the theorem using  $\frac{1}{\sqrt{2\lambda}}\|U_{\lambda,\perp}B\|_F$ , where  $U_{\lambda,\perp}$  denotes the bottom  $m_S-p$  left singular vectors of a sketch SA of A, where p is such that  $\lambda$  is between  $\sigma_{p+1}^2(SA)$  and  $\sigma_p^2(SA)$ . (Here we use  $p=\mathrm{rank}A$  when  $\lambda \leq \sigma_{\mathrm{rank}A}^2$ .) We could write this roughly as  $\|A_{-p}A_{-p}^+B\|/\|A_{-p}\|$ , where  $A_{-p}$  denotes A minus its best rank p approximation. It is the part of B we are "giving up" by including a ridge term. Proofs for this section are deferred to Appendix B.

**Theorem 1.4** (Dynamic Ridge Regression). Given an  $n \times d$  matrix A of rank k, an  $n \times d'$  matrix B, error parameter  $\varepsilon > 0$ , and ridge parameter  $\lambda$ , let  $\kappa_{\lambda}^2 = (\lambda + \sigma_1^2(A))/(\lambda + \sigma_k^2(A))$  be the ridge condition number of A, and let  $\psi_{\lambda} = \|A\|_F^2/(\lambda + \sigma_k^2(A))$ . Further, let  $X^*$  be the optimal ridge regression solution, i.e.,  $X^* = \operatorname{argmin}_X \|AX - B\|_F^2 + \lambda \|X\|_F^2$ .

Then there is a data structure supporting turnstile updates of the entries of A in  $O(\log(nd))$  time, and an algorithm using that data structure that computes a sample SA of  $m = O\left(\kappa_{\lambda}^2 \psi_{\lambda} \log(nd)/\varepsilon^2\right)$  rows of A, where  $S \in \mathbb{R}^{m \times n}$  is a sampling matrix, and outputs  $\tilde{X} \in \mathbb{R}^{m \times d'}$ , such that with probability 99/100,

$$\|A^{\top}S^{\top}\tilde{X} - X^*\|_F \le \varepsilon \left(1 + 2\gamma\right) \|X^*\|_F + \frac{\varepsilon}{\sqrt{\lambda}} \|U_{k,\perp}B\|_F,$$

where  $U_{\lambda,\perp}B$  is the projection of B onto the subspace corresponding to the singular values of SA less than  $\sqrt{\lambda}$ ; and  $\gamma^2 = \frac{\|B\|_F^2}{\|AA^+X^*\|_F^2}$  is a problem dependent parameter.

Further, the running time of the algorithm is  $\tilde{O}\left(d'\varepsilon^{-4}\psi_{\lambda}^2\kappa_{\lambda}^2\log(d)\right)$ . Finally, for all  $i\in[d]$ , and  $j\in[d']$ , an entry  $(A^{\top}S^{\top}\tilde{X})_{i,j}$  can be computed in  $O(m\log(nd))$  time.

We note that the "numerical" quantities  $\kappa_{\lambda}$  and  $\psi_{\lambda}$  are

decreasing in  $\lambda$ . When  $\lambda$  is within a constant factor of  $\|A\|^2$ ,  $\psi_{\lambda}$  is within a constant factor of the *stable rank*  $\|A\|_F^2/\|A\|^2$ , where the stable rank is always at most rank(A). We also note that in the theorem, and the remainder of the paper, a *row sampling matrix* S has rows that are multiples of natural basis vectors, so that SA is a (weighted) sample of the rows of A. A column sampling matrix is defined similarly.

Concurrent Work on Ridge Regression. In an independent and concurrent work, Gilyén, Song and Tang (Gilyén et al., 2020) obtain a roughly comparable classical algorithm for regression, assuming access to the tree data structure, which runs in time  $\tilde{O}\left(\frac{\|A\|_F^6\|A\|_2^2}{\|A^+\|_2^8\varepsilon^4}\right)$ , or in the notation above,  $\tilde{O}(\varepsilon^{-4}\psi_\lambda^3\kappa^2)$ , for the special case of d'=1. Their algorithm is based on Stochastic Gradient Descent.

Next, we describe our results for low-rank approximation. We obtain a dynamic algorithm (Algorithm 3) for approximating A with a rank-k matrix, for a given k, and a data structure for sampling from it, in the vein of (Tang, 2019). At a high-level, as with our ridge regression algorithm, we first sample rows proportional to their squared Euclidean norm (length-squared sampling) and then sample a subset of columns resulting in a small submatrix with  $\tilde{O}(\varepsilon^{-2}k)$  rows and columns. We then compute the SVD of this matrix, and then work back up to A with more sampling and a QR factorization. The key component in our algorithm and analysis is using  $Projection-Cost\ Preserving$  sketches (see Definition A.9). These enable us to preserve the Frobenius cost of projections onto all rank-k subspaces simultaneously. As a result, we obtain the following theorem:

**Theorem 1.5** (Sampling from a low-rank approximation). Given an  $n \times d$  matrix A for which a sampling data structure has been maintained, target rank  $k \in [d]$  and error parameter  $\varepsilon > 0$ , we can find sampling matrices S and R, and rank-k matrix W, such that  $\|ARWSA - A\|_F \le (1 + O(\varepsilon))\|A - A_k\|_F$ . Further, the running time is  $\tilde{O}(\varepsilon^{-6}k^3 + \varepsilon^{-4}\psi_\lambda(\psi_\lambda + k^2 + k\psi_k))$ , where  $\psi_\lambda$  is as in Theorem 1.4, and  $\psi_k = \frac{\|A\|_F^2}{\sigma_k(A)^2}$ . Given  $j \in [d]$ , a random index  $i \in [n]$  with probability distribution  $(ARWSA)_{ij}^2/\|ARWSA)_{*,j}\|^2$  can be generated in expected time  $\tilde{O}(\psi_k + k^2\varepsilon^{-2}\kappa^2)$ , where  $\kappa = \sigma_1(A) \cdot \sigma_{\mathrm{rank}(A)}(A)$ .

Here if the assumption  $||A_k||_F^2 \ge \varepsilon ||A||_F^2$  does not hold, the trivial solution 0 satisfies the relative error target and we assume the resulting approximation is not worth sampling:

$$||A - 0||_F^2 \le \frac{1}{1 - \varepsilon} (||A||_F^2 - ||A_k||_F^2) = \frac{1}{1 - \varepsilon} ||A - A_k||_F^2.$$

This result is directly comparable to Tang's algorithm (Tang, 2019) for recommender systems which again needs query

time that is a large polynomial in  $k, \kappa$  and  $\varepsilon^{-1}$ . Our algorithm returns a relative error approximation, a rank-k approximation within  $1+\varepsilon$  of the best rank-k approximation; Tang's algorithm has additive error, with a bound more like  $\|A-A_k\|_F + \varepsilon \|A\|_F$ . Finally, we note that  $\psi_k \leq \kappa^2$  and for several settings of k can be significantly smaller.

For ease of comparison we summarize our results in Table 1.

#### 1.2. Related Work

**Matrix Sketching.** The *sketch and solve* paradigm (Clarkson & Woodruff, 2015; Woodruff, 2014) was designed to reduce the dimensionality of a problem, while maintaining enough structure such that a solution to the smaller problem remains an approximate solution the original one. This approach has been pivotal in speeding up basic linear algebra primitives such as least-squares regression (Sarlós, 2006; Rokhlin & Tygert, 2008; Clarkson & Woodruff, 2015),  $\ell_p$ regression (Cohen & Peng, 2015; Wang & Woodruff, 2019), low-rank approximation (Nelson & Nguyên, 2013; Cohen et al., 2017; Li & Woodruff, 2020), linear and semi-definite programming (Cohen et al., 2019; Jiang et al., 2020b;a) and solving non-convex optimization problems such as  $\ell_n$ low-rank approximation (Song et al., 2017; 2019; Ban et al., 2019) and training neural networks (Bakshi et al., 2019b; Brand et al., 2020). For a comprehensive overview we refer the reader to the aforementioned papers and citations therein. Several applications use rank computation, finding a full rank subset of rows/columns, leverage score sampling, and computing subspace embeddings as key algorithmic primitives.

Sublinear Algorithms and Quantum Linear Algebra. Recently, there has been a flurry of work on sublinear time algorithms for structured linear algebra problems (Musco & Woodruff, 2017; Shi & Woodruff, 2019; Balcan et al., 2019; Bakshi et al., 2020) and quantum linear algebra (Harrow et al., 2009; Gilyén et al., 2019; Lloyd et al., 2014; Kerenidis & Prakash, 2016; Dunjko & Wittek, 2020). The unifying goal of these works is to avoid reading the entire input to solve tasks such as linear system solving, regression and low-rank approximation. The work on sublinear algorithms assumes the input is drawn from special classes of matrices, such as positive semi-definite matrices (Musco & Woodruff, 2017; Bakshi et al., 2019a), distance matrices (Bakshi & Woodruff, 2018; Indyk et al., 2019) and Toeplitz matrices (Lawrence et al., 2020), whereas the quantum algorithms (and their de-quantized analogues) assume access to data structures that admit efficient sampling (Tang, 2019; Gilyén et al., 2018; Chia et al., 2020).

The work of Gilyén, Lloyd and Tang (Gilyén et al., 2018) on low-rank least squares produces a data structure as output:

given index  $i\in[d]=\{1,\ldots,d\}$ , the data structure returns entry  $x_i'$  of  $x'\in\mathbb{R}^d$ , which is an approximation to the solution  $x^*$  of  $\min_{x\in\mathbb{R}^d}\|Ax-b\|$ , where  $b\in\mathbb{R}^n$ . The error bound is  $\|x'-x^*\|\leq \varepsilon \|x^*\|$ , for given  $\varepsilon>0$ . This requires the condition that  $\|Ax^*-b\|/\|Ax^*\|$  is bounded above by a constant. Subsequent work (Chia et al., 2020) removes this requirement, and both results obtain data structures that need space polynomial in  $\mathrm{rank}(A),\ \varepsilon,\ \kappa(A),^1$  and other parameters.

The work (Tang, 2019) also produces a data structure, that supports sampling relevant to the setting of recommender systems: the nonzero entries of the input matrix A are a subset of the entries of a matrix P of, for example, user preferences. An entry  $A_{ij} \in [0,1]$  is one if user j strongly prefers product i, and zero if user j definitely does not like product i. It is assumed that P is well-approximated by a matrix of some small rank k. The goal is to estimate Pusing A; one way to make that estimate effective, without simply returning all entries of P, is to create a data structure so that given j, a random index i is returned, where i is returned with probability  $\hat{a}_{ij}^2/\|\hat{A}_{*,j}\|^2$ . Here  $\hat{A}_{*,j}$  is the j'th column of  $\hat{A}$  (and  $\hat{a}_{ij}$  an entry), where  $\hat{A}$  is a good rank-k approximation to A, and therefore, under appropriate assumptions, to P. The estimate  $\hat{A}$  is regarded as a good approximation if  $\|\hat{A} - A\|_F \le (1+\varepsilon)\|A - [A]_k\|_F$ , where  $[A]_k$  is the matrix of rank k closest to A in Frobenius norm. Here  $\varepsilon$  is a given error parameter. As shown in (Tang, 2019), this condition (or indeed, a weaker one) implies that the described sampler is useful in the context of recommender systems.

#### 2. Outline

The next section gives some notation and mathematical preliminaries, in particular regarding leverage-score and length-squared sampling. This is followed by descriptions of our data structures and algorithms, and then by our computational experiments. The appendices give some extensive descriptions, proofs of theorems, and in Appendix D, some additional experiments.

#### 3. Preliminaries

Let  $X^+$  denote the Moore-Penrose pseudo-inverse of matrix X, equal to  $V\Sigma^{-1}U^\top$  when X has thin SVD  $X=U\Sigma V^\top$ , so that  $\Sigma$  is a square invertible matrix. We note that  $X^+=(X^\top X)^+X^\top=X^\top(XX^\top)^+$  and  $X^+XX^\top=X^\top$ , which is provable using the SVDs of X and  $X^+$ . Also, if X has full column rank, so that Y is square, then  $X^+$  is a left inverse of X, that is,  $X^+X=I_d$ , where d is the

<sup>&</sup>lt;sup>1</sup>Throughout, we define  $\kappa(A) = ||A|| ||A^+||$ , that is, the ratio of largest to smallest *nonzero* singular values of A, so that, in particular, it will never be infinite or undefined.

Table 1. Comparison of our results and prior work. Let the target error be  $\varepsilon$ , target rank be k and let  $\psi_k = \|A\|_F^2/\sigma_k(A)^2$ , where  $\sigma_k$  is the k-th singular value of the input matrix. Also,  $\hat{\sigma}_k \leq 1/\|A^+\|$ ,  $\hat{\sigma}_1 \geq \|A\|$ , d' is the number of columns of B for multiple-response, and  $\eta$  denotes some numerical properties of A. To avoid numerous parameters, we state our results by setting  $\lambda = \Theta(\|A\|_2^2)$  in the corresponding theorems.

Problem		Time	Prior Work		
	Update	Query	Update	Query	
Ridge Regression	$O(\log(n))$	$\tilde{O}\left(\frac{d'\kappa^3 \ A\ _F^2 \log(d)}{\varepsilon^4 \ A\ _2^2}\right)$	$O(\log(n))$	$\tilde{O}\left(\frac{k^6 \ A\ _F^6 \kappa^{16}}{\varepsilon^6}\right)$	
		Thm. 1.4		(Gilyén et al., 2018) $\tilde{O}\left(\frac{\ A\ _F^8 \kappa(A)^2}{(\sigma_{\min}^6 \varepsilon^4)}\right)$ (Gilyén et al., 2020)	
Low Rank Sampling	$O(\log(n))$	$\tilde{O}\left(\frac{\ A\ _F^2 \left(\frac{\ A\ _F^2}{\ A\ _2^2} + k^2 + k\psi_k\right)}{\varepsilon^4 \ A\ _2^2} + \frac{k^3}{\varepsilon^6}\right)$	$O(\log(n))$	$\Omega(\operatorname{poly}(\kappa k \varepsilon^{-1} \eta))$	
		Thm. 1.5		(Tang, 2019)	

number of columns of X. Let  $\|X\|$  denote the spectral (operator) norm of X. Let  $\kappa(X) = \|X^+\| \|X\|$  denote the condition number of X. We write  $a \pm b$  to denote the set  $\{c \mid |c-a| \leq |b|\}$ , and  $c=a \pm b$  to denote the condition that c is in the set  $a \pm b$ . Let  $[m] = \{1, 2, \ldots, m\}$  for an integer m.

As mentioned,  $\operatorname{nnz}(A)$  is the number of nonzero entries of A, and we assume  $\operatorname{nnz}(A) \geq n$ , which can be ensured by removing any rows of A that only contain zeros. We let  $[A]_k$  or sometimes  $A_k$  denote the best rank-k approximation to A. Let  $0_{a \times b} \in \mathbb{R}^{a \times b}$  have all entries equal to zero, and similarly  $0_a \in \mathbb{R}^a$  denotes the zero vector. Further, for an  $n \times d$  matrix A and a subset S of [n], we use the notation  $A_{|S|}$  to denote the restriction of the rows of A to the subset indexed by S. As mentioned,  $n^\omega$  is the time needed to multiply two  $n \times n$  matrices.

**Lemma 3.1** (Oblivious Subspace Embedding Theorem 7.4 (Chepurko et al., 2022)). For given matrix  $A \in \mathbb{R}^{n \times d}$  with  $k = \operatorname{rank}(A)$ , there exists an oblivious sketching matrix S that samples  $m = O(\varepsilon_0^{-2}k\log k)$  rows of A such that with probability at least 99/100, for all  $x \in \mathbb{R}^d$ , S is an  $\varepsilon_0$ -subspace embedding, that is,  $||SAx|| = (1 \pm \varepsilon_0)||Ax||$ . Further, the matrix SA can be computed in  $O(\operatorname{nnz}(A) + k^{\omega}\operatorname{poly}(\log\log(k)) + \operatorname{poly}(1/\varepsilon_0)k^{2+o(1)})$  time.

We obtain the following data structure for leverage-score sampling. We provide a statement of its properties below, but defer the description of the algorithm and proof to the supplementary material. While leverage-score sampling is well-known, we give an algorithm for completeness; also, our algorithm removes a log factor in some terms in the runtime, due to our use of the sketch of Lemma 3.1.

**Theorem 3.2** (Leverage Score Data Structure). Let k = rank(A), and choose  $\mu_s \geq 1$ . Then, Algorithm 5 (LEVSAMPLE $(A, \mu_s, v)$ ) uses space  $O(n+k^{\omega} \log \log(nd))$ ,

not counting the space to store A, and runs in time

$$O(\mu_s \operatorname{nnz}(A) + k^{\omega} \operatorname{poly}(\log \log(k)) + k^{2+o(1)} + vkn^{1/\mu_s}),$$

and outputs a leverage score sketching matrix L, which samples v rows of A with probability proportional to their leverage scores. (It also outputs a column selector  $\Lambda$ , selecting an orthogonal basis of the column space of A.)

**Definition 3.3** (Ridge Leverage-score Sample, Statistical Dimension). Let A be such that  $k = \operatorname{rank}(A)$ , and suppose A has thin SVD  $A = U\Sigma V^{\top}$ , implying  $\Sigma \in \mathbb{R}^{k\times k}$ . For  $\lambda > 0$ , let  $A_{(\lambda)} = \begin{bmatrix} A \\ \sqrt{\lambda}VV^{\top} \end{bmatrix}$  and  $A_{(\lambda)}$  has SVD  $A_{(\lambda)} = \begin{bmatrix} U\Sigma D \\ \sqrt{\lambda}VD \end{bmatrix} D^{-1}V^{\top}$ , where  $D = (\Sigma^2 + \lambda I_k)^{-1/2}$ . Call  $S \subset [n]$  a ridge leverage-score sample of A if each  $i \in S$  is chosen independently with probability at least  $\|U_{i,*}\Sigma D\|^2/\operatorname{sd}_{\lambda}(A)$ , where the statistical dimension  $\operatorname{sd}_{\lambda}(A) = \|U\Sigma D\|_F^2 = \sum_{i \in [d]} \sigma_i^2/(\lambda + \sigma_i^2)$ , recalling that  $U\Sigma D$  comprises the top n rows of the left singular matrix of  $A_{(\lambda)}$ .

We can also use *length-squared sampling* to obtain subspace embeddings. In Section 4 we will give a data structure and algorithm that implements length-squared sampling. We defer the analysis to Appendix A.

**Definition 3.4** (Length-squared sample). Let  $A \in \mathbb{R}^{n \times d}$ ,  $\lambda \geq 0$ , and  $A_{(\lambda)}$  be as in Lemma A.5. For given m, let matrix  $L \in \mathbb{R}^{m \times n}$  be chosen by picking each row of L to be  $e_i^{\top}/\sqrt{p_i m}$ , where  $e_i \in \mathbb{R}^n$  is the i'th standard basis vector, and picking  $i \in [n]$  with probability  $p_i \leftarrow \|A_{i,*}\|^2/\|A\|_F^2$ .

We obtain the corresponding lemma for length-squared sampling and defer the proof to Appendix A.

**Lemma 3.5** (Length-squared sketch). Given a matrix  $A \in \mathbb{R}^{n \times d}$  and a sample size parameter  $v \in [n]$ , let  $m = O\left(v\|A_{(\lambda)}^+\|^2\|A\|_F^2/\operatorname{sd}_{\lambda}(A)\right)$ . Then, with probability at

least 99/100, the set of m length-squared samples contains a ridge leverage-score sample of A of size v.

# 4. Dynamic Data Structures for Ridge Regression

In this section, we describe our dynamic data structures, and then our algorithm for solving Ridge Regression problems. Given an input matrix  $A \in \mathbb{R}^{n \times d}$ , our data structure can be maintained under insertions and deletions (and changes) in  $O(\log(nd))$  time, such that sampling a row or column with probability proportional to its squared length can be done in  $O(\log(nd))$  time. The data structure is used for solving both ridge regression and LRA (Low-Rank Approximation) problems.

First, we start with a simple folklore data structure.

**Lemma 4.1.** Given  $\ell$  real values  $\{u_i\}_{i\in[\ell]}$ , there is a data structure using storage  $O(\ell)$ , so that  $L = \sum_{i\in[\ell]} u_i^2$  can be maintained, and such that a random i can be chosen with probability  $u_i^2/L$  in time  $O(\log \ell)$ . Values can be inserted, deleted, or changed in the data structure in time  $O(\log \ell)$ .

The implementation of this data structure is discussed in Appendix B. We use it in our data structure DYNSAMP(A), given below, which is used in LENSQSAMPLE, Alg. 1, to sample rows and columns of A.

**Definition 4.2.** DYNSAMP(A) is a data structure that, for  $A \in \mathbb{R}^{n \times d}$ , comprises:

- For each row of A, the data structure of Lemma 4.1 for the nonzero entries of the row or column.
- For the rows of A, the data structure of Lemma 4.1 for their lengths.
- For given i, j, a data structure supporting access to the value of entry a<sub>ij</sub> of A in O(1) time.

```
Algorithm 1 LENSQSAMPLE(DS, S = \text{null}, m_S, m_R)
```

Input: DS = DYNSAMP(A) (Def. 4.2) for  $A \in \mathbb{R}^{n \times d}$ , sample sizes  $m_S$ ,  $m_R$ 

**Output:** Sampling matrices  $S \in \mathbb{R}^{m_S \times n}$ ,  $R \in \mathbb{R}^{d \times m_R}$ 

- 1: if  $S == \mathbf{null}$ 
  - Use DS to build row sampler  $S \in \mathbb{R}^{m_S \times n}$  of A
- 2: Use DS and S to build column sampler  $\mathbb{R}^{d \times m_R}$  of SA {cf. Lemma 4.3}
- 3: return S, R

**Lemma 4.3.** DYNSAMP(A) can be maintained under turnstile updates of A in  $O(\log(nd))$  time. Using DYNSAMP(A), rows can be chosen at random with row  $i \in [n]$  chosen with probability  $\|A_{i,*}\|^2/\|A\|_F^2$  in  $O(\log(nd))$  time.

```
Algorithm 2 RIDGEREGDYN(DS, B, \hat{\sigma}_k, \hat{\sigma}_1, \varepsilon, \lambda)
```

**Input:** DS = DYNSAMP(A),  $B \in \mathbb{R}^{n \times d'}$ ,  $\hat{\sigma}_k \leq 1/\|A^+\|$ ,  $\hat{\sigma}_1 \geq \|A\|$ ,  $\varepsilon$  an error parameter,  $\lambda$  a ridge weight

**Output:** Data for approximate ridge regression solution  $A^{\top}S^{\top}\tilde{X}$  where S is a sampling matrix

```
1: Z_{\lambda} \leftarrow 1/\sqrt{\lambda + \hat{\sigma}_{k}^{2}}, \hat{\kappa} \leftarrow Z_{\lambda}\sqrt{\lambda + \hat{\sigma}_{1}^{2}}

2: Choose m_{S} = O(\varepsilon^{-2}\hat{\kappa}^{2}Z_{\lambda}^{2}\|A\|_{F}^{2}\log(d)), m_{R} = O(\hat{m}_{R}Z_{\lambda}^{2}\|A\|_{F}^{2}), where \hat{m}_{R} = O(\varepsilon^{-2}\log m_{S})

3: S, R \leftarrow \text{Lensqsample(DS, null}, m_{S}, m_{R}) {cf. Alg. 1;}

4: \tilde{X} \leftarrow (SARR^{\top}A^{\top}S^{\top} + \lambda I_{m_{S}})^{-1}SB {Solve using conjugate gradient}

5: return \tilde{X}, S {approximate ridge regression solution is A^{\top}S^{\top}\tilde{X}}
```

If  $S \in \mathbb{R}^{m \times n}$  is a sampling matrix, so that SA has rows that are each a multiple of a row of A, then c columns can be sampled from SA using DYNSAMP(A) in  $O((c+m)\log(nd))$  time, with the column  $j \in [d]$  chosen with probability  $\|(SA)_{*,j}\|^2/\|SA\|_E^2$ .

We designate the algorithm of Lemma 3.5 as LenSqSample, as given at a high level in Algorithm 1, and in more detail in the proof of Lemma 4.3 in Appendix B.

This simple data structure and sampling scheme will be used to solve ridge regression problems, via Algorithm 2. Its analysis, which proves Theorem 1.4, is given in Appendix B.

# 5. Sampling from a Low-Rank Approximation

Our algorithm for low-rank approximation is BUILDLOWRANKFACTORS, Algorithm 3, given below. As discussed in the introduction, it uses LENSQSAMPLE, Algorithm 1, to reduce to a matrix whose size is independent of the input size, beyond log factors, as well as Projection-Cost Preserving sketches, QR factorization, and leverage-score sampling. Its analysis, proving Theorem 1.5, is given in Appendix C.

#### 6. Experiments

We evaluate the empirical performance of our algorithm on both synthetic and real-world datasets. All of our experiments were done in Python and conducted on a laptop with a 1.90GHz CPU and 16GB RAM. Prior work (Arrazola et al., 2020) suggests the tree data structure is only faster than the built-in sampling function when the matrix size  $\max\{n,d\}$  is larger than  $10^6$ . Hence we follow the implementation in (Arrazola et al., 2020) that directly uses the built-in function. For a fair comparison, we also modified the code in (Arrazola et al., 2020), which reduces the time

### Algorithm 3 BUILDLOWRANKFACTORS (DYNSAMPLER, $k, \hat{\sigma}_k, \hat{\sigma}_k, \varepsilon, \tau$ )

Input: DYNSAMPLER = DYNSAMP(A) (Def. 4.2) for  $A \in \mathbb{R}^{n \times d}$ , k target rank,  $\hat{\sigma}_k \leq 1/\|A^+\|$ ,  $\hat{\sigma}_k \leq \sigma_k(A)$ ,  $\varepsilon$  an error parameter,  $\tau$  estimate of  $\|A - A_k\|_F^2$ , where  $A_k$  is the best rank- $\bar{k}$  approximation to A

Output: Small matrix W and sampling matrices S,Rso that rank(ARWSA) = k and  $||ARWSA - A|| \le (1 + \varepsilon)||A - A_k||$ 

- 1:  $\lambda \leftarrow \tau/k, Z_{\lambda} \leftarrow 1/\sqrt{\lambda + \hat{\sigma}_{k}^{2}}, Z_{k} \leftarrow 1/\hat{\sigma}_{k}$ 2: Choose  $m_R = m_S = O(\hat{m}_S Z_{\lambda}^2 ||A||_F^2)$ ,
- where  $\hat{m}_S = O(\varepsilon^{-2} \log k)$
- 3:  $S, R_1 \leftarrow \text{LenSqSample}(\text{DynSampler}, \mathbf{null}, m_S, m_R)$
- 4: Apply Alg. 1 and Thm. 1 of (Cohen et al., 2017) to  $SAR_1$ , get col. sampler  $R_2$  {  $m_{R_2} = O(\varepsilon^{-2}k \log k)$ } 5: Apply Alg. 1 and Thm. 1 of (Cohen et al., 2017) to  $SAR_1R_2$ ,
- get row sampler  $S_2$  {  $m_{S_2} = O(\varepsilon^{-2}k \log k)$  }
- 6:  $V \leftarrow \text{top-}k \text{ right singular matrix of } S_2S\bar{A}R_1R_2$
- 7: U,  $\leftarrow QR(SAR_1R_2V)$

 $\{U \text{ has orthonormal cols}, SAR_1R_2V = UC \text{ for matrix } C\}$ 

- 8: Choose  $m_{R_3} = O(\hat{m}_{R_3} \varepsilon^{-1} Z_k^2 \|A\|_F^2)$ , where  $\hat{m}_{R_3} = O(\varepsilon_0^{-2} \log k + \varepsilon^{-1})$ ,  $\varepsilon_0$  a small constant 9:  $R_3 \leftarrow \text{LenSqSample}(\text{DynSampler}, S, m_S, m_{R_3})$
- 10: Let f(k, C) be the function returning the value
- $$\begin{split} m_{R_4} &= O(\varepsilon_0^{-2} k \log k + \varepsilon^{-1} k) \\ \text{11: } R_4^\top, &\quad \leftarrow \text{LevSample}((U^\top SAR_3)^\top, \log(m_{R_3}), f()) \end{split}$$
  {Alg. 5}
- 12:  $R \leftarrow R_3 R_4$
- 13:  $W \leftarrow (U^{\top}SAR)^+U^{\top}$
- 14: return  $\hat{W}$ , S, R

to maintain the data structure by roughly 30x. For each experiment, we took an average over 10 independent trials.

We note that we do not compare with classical sketching algorithms for several reasons. First, there is no classical contender with the same functionality as ours. This is because our dynamic algorithms support operations not seen elsewhere: sublinear work for regression and low-rank approximation, using simple fast data structures that allow, as special cases, row-wise or column-wise updates. Second, unlike dynamic algorithms where a sketch is maintained, our algorithms are not vulnerable to updates based on prior outputs, whether adversarially, or due to use in the inner loop of an optimization problem. This is because our algorithms are based on independent sampling from the exact input matrix.

#### 6.1. Low-Rank Approximation

We conduct experiments on the following datasets:

• KOS data.<sup>2</sup> A word frequency dataset. The matrix represents word frequencies in blogs and has dimensions

- $3430 \times 6906$  with 353160 non-zero entries.
- MovieLens 100K. (Harper & Konstan, 2016) A movie ratings dataset, which consists of a preference matrix with 100,000 ratings from 611 users across 9,724 movies.

We compare our algorithms with the implementations in (Arrazola et al., 2020), which are based on the algorithms in (Frieze et al., 2004) and (Tang, 2019). We refer to this algorithm as ADBL henceforth. For the KOS dataset, we set the number of sampled rows and columns to be (r,c) = (500,700) for both algorithms. For the MovieLens dataset we set (r, c) = (300, 500). We define the error  $\varepsilon =$  $||A - Y||_F / ||A - A_k||_F - 1$ , where Y is the algorithm's output and  $A_k$  is the best k-rank approximation. Since the regime of interest is  $k \ll n$ , we vary k among  $\{10, 15, 20\}$ .

The results are shown in Table 2. We first report the total runtime, which includes the time to maintain the data structure and then compute the low-rank approximation. We also report the query time, which excludes the time to maintain the data structure. From the table we see that both algorithms can achieve  $\varepsilon \approx 0.05$  in all cases. The query time of ours is about 6x-faster than the ADBL algorithm in (Arrazola et al., 2020), and even for the total time, our algorithm is much faster than the SVD. Although the accuracy of ours is slightly worse, in Appendix D.1 we show by increasing the sample size slightly, our algorithm achieves the same accuracy as ADBL (Arrazola et al., 2020), but still has a faster runtime.

We remark that the reason our algorithm only needs half of the time to compute the sampling probabilities is that we only need to sample rows or columns according to their squared length, but the algorithm in (Arrazola et al., 2020) also needs to sample entries for each sampled row according to the squared values of the entries.

#### 6.2. Ridge Regression

In this section, we consider the problem

$$X^* := \min_{X \in \mathbb{D}^{d \times d'}} \|AX - B\|_F^2 + \lambda \|X\|_F^2,$$

where  $A \in \mathbb{R}^{n \times d}$ ,  $B \in \mathbb{R}^{n \times d'}$ . We do experiments on the following dataset with  $\lambda = 1$ :

- Synthetic data. We generate the rank-k matrix Aas (Arrazola et al., 2020) do. Particularly, suppose the SVD of A is  $A = U\Sigma V^{\top}$ . We first sample an  $n \times k$ Gaussian matrix, then we perform a QR-decomposition G = QR, where Q is an  $n \times k$  orthogonal matrix. We then simply set U = Q and then use a similar way to generate V. We set  $A \in \mathbb{R}^{7000 \times 9000}$ ,  $B \in \mathbb{R}^{7000 \times 1}$ .
- **YearPrediction.**<sup>3</sup> A dataset that collects 515345 songs

<sup>&</sup>lt;sup>2</sup>The Bag of Words Data Set from the UCI Machine Learning Repository.

<sup>&</sup>lt;sup>3</sup>YearPredictionMSD Data Set

Table 2. Performance of our algorithm and ADBL on MovieLen 100K and KOS data, respectively

and KOS data, respectively.					
	k = 10	k = 15	k = 20		
$\varepsilon$ (Ours)	0.0416	0.0557	0.0653		
$\varepsilon(ADBL)$	0.0262	0.0424	0.0538		
Runtime	0.125s	0.131s	0.135s		
(Ours, Query)	0.1238	0.1318			
Runtime	0.181s	0.183s	0.184s		
(Ours, Total)	0.1013	0.1033	0.10-3		
Runtime	0.867s	0.913s	1.024s		
(ADBL, Query)	0.0075	0.5135			
Runtime	0.968s	1.003s	1.099s		
(ADBL, Total)	0.7000		1.0555		
Runtime of SVD	2.500s				
	k = 10	k = 15	k = 20		
$\varepsilon$ (Ours)	0.0397	0.0478	0.0581		
$\varepsilon(ADBL)$	0.0186	0.0295	0.0350		
	0.0100	0.0273	0.0550		
Runtime					
Runtime (Ours, Query)	0.292s	0.296s	0.295s		
	0.292s	0.296s	0.295s		
(Ours, Query)					
(Ours, Query) Runtime (Ours, Total) Runtime	0.292s 0.452s	0.296s 0.455s	0.295s 0.452s		
(Ours, Query)  Runtime (Ours, Total)  Runtime (ADBL, Query)	0.292s	0.296s	0.295s		
(Ours, Query)  Runtime (Ours, Total)  Runtime (ADBL, Query)  Runtime	0.292s 0.452s 1.501s	0.296s 0.455s 1.643s	0.295s 0.452s 1.580s		
(Ours, Query)  Runtime (Ours, Total)  Runtime (ADBL, Query)	0.292s 0.452s	0.296s 0.455s	0.295s 0.452s		

and each song has 90 attributes. The task here is to predict the release year of the song.  $A \in \mathbb{R}^{515345 \times 90}$ ,  $B \in \mathbb{R}^{515345 \times 1}$ .

• **PEMS data.** <sup>4</sup> The data describes the occupancy rate of different car lanes of San Francisco bay area freeways. Each row is the time series for a single day. The task on this dataset is to classify each observed day as the correct day of the week, from Monday to Sunday.  $A \in \mathbb{R}^{440 \times 138672}, B \in \mathbb{R}^{440 \times 1}$ .

We define the error  $\varepsilon = \|X - X^*\|_F / \|X^*\|_F$ , given the algorithm output X. For synthetic data, we set the number of sampled rows and columns to be r and c. For the YearPrediction data, the number of columns is small, and hence we only do row sampling, and likewise, for the PEMS data, we only do column sampling. We did not find an implementation for the ridge regression problem in the previous related work. Therefore, here we list the time to compute the closed-form optimal solution  $X^* = (A^T A + \lambda I)^{-1} A^T B$  or  $X^* = A^T (AA^T + \lambda I)^{-1} B$ , as a reference.

The results are shown in Table 3 and 4. From the tables we can see that for synthetic data, the algorithm can achieve an error  $\varepsilon < 0.1$  when only sampling less than 10% of the rows

*Table 3.* Performance of our algorithm on synthetic data.

(r,c)	300,500	500,800	1000, 1500	
$\varepsilon$ (Ours)	0.1392	0.0953	0.0792	
Runtime (Query)	0.021s	0.042s	0.148s	
Runtime (Total)	0.557s	0.568s	0.667s	
Exact X*		24.074s		

and columns. Also, the total runtime is about 40x-faster than computing the exact solution. For the YearPrediction and PEMS data, the bottleneck of the algorithm becomes the time to compute the sample probabilities, but the query time is still very fast and we can achieve an error  $\varepsilon < 0.1$  when only sampling a small fraction of the rows or columns.

Table 4. Performance of our algorithm on YearPrediction data and PEMS data, respectively.

ita, respectively.				
r =	1000	3000	5000	
$\varepsilon$ (Ours)	0.1070	0.0633	0.0447	
Runtime	0.031s	0.037s	0.059s	
(Query)	0.0518			
Runtime	0.213s	0.229s	0.245s	
(Total)	0.2138	0.2298	0.2438	
Exact X*	0.251s			
c =	15000	25000	35000	
$\varepsilon$ (Ours)	0.1778	0.1397	0.1130	
Runtime	0.234s	0.381s	0.532s	
(Query)		0.3618	0.5528	
Runtime	0.473s	0.628s	0.777s	
(Total)	0.4738	0.0268		
Exact X*	0.972s			

#### Acknowledgements.

Honghao Lin and David Woodruff would like to thank for partial support from the National Science Foundation (NSF) under Grant No. CCF-1815840.

#### References

Arrazola, J. M., Delgado, A., Bardhan, B. R., and Lloyd, S. Quantum-inspired algorithms in practice. *Quantum*, 4:307, 2020. URL https://doi.org/10.22331/q-2020-08-13-307.

Avron, H., Clarkson, K. L., and Woodruff, D. P. Sharper bounds for regularized data fitting. In *RAN-DOM '17: 21st International Workshop on Ran-domization and Computation*, 2017. URL https:

<sup>&</sup>lt;sup>4</sup>PEMS-SF Data Set

- //arxiv.org/abs/1611.03225. Full version at https://arxiv.org/abs/1611.03225.
- Bakshi, A. and Woodruff, D. Sublinear time low-rank approximation of distance matrices. In *Advances in Neural Information Processing Systems*, pp. 3782–3792, 2018.
- Bakshi, A., Chepurko, N., and Woodruff, D. P. Robust and sample optimal algorithms for psd low-rank approximation. *arXiv preprint arXiv:1912.04177*, 2019a.
- Bakshi, A., Jayaram, R., and Woodruff, D. P. Learning two layer rectified neural networks in polynomial time. In *Conference on Learning Theory*, pp. 195–268. PMLR, 2019b.
- Bakshi, A., Chepurko, N., and Jayaram, R. Testing positive semi-definiteness via random submatrices. *arXiv* preprint *arXiv*:2005.06441, 2020.
- Balcan, M.-F., Li, Y., Woodruff, D. P., and Zhang, H. Testing matrix rank, optimally. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 727–746. SIAM, 2019.
- Ban, F., Bhattiprolu, V., Bringmann, K., Kolev, P., Lee, E., and Woodruff, D. P. A PTAS for  $\ell_p$ -low rank approximation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 747–766. SIAM, 2019.
- Berry, D. W., Childs, A. M., and Kothari, R. Hamiltonian simulation with nearly optimal dependence on all parameters. In 2015 IEEE 56th Annual Symposium on Foundations of Computer Science, pp. 792–809, 2015.
- Boutsidis, C., Woodruff, D. P., and Zhong, P. Optimal principal component analysis in distributed and streaming models. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pp. 236–249, 2016.
- Brand, J. v. d., Peng, B., Song, Z., and Weinstein, O. Training (overparametrized) neural networks in near-linear time. *arXiv* preprint arXiv:2006.11648, 2020.
- Brandão, F. G. S. L., Kalev, A., Li, T., Lin, C. Y.-Y., Svore, K. M., and Wu, X. Quantum SDP Solvers: Large Speed-Ups, Optimality, and Applications to Quantum Learning. In Baier, C., Chatzigiannakis, I., Flocchini, P., and Leonardi, S. (eds.), 46th International Colloquium on Automata, Languages, and Programming (ICALP 2019), volume 132 of Leibniz International Proceedings in Informatics (LIPIcs), pp. 27:1–27:14, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-109-2. doi: 10.4230/LIPIcs.ICALP.2019.27. URL http://drops.dagstuhl.de/opus/volltexte/2019/10603.

- Chepurko, N., Clarkson, K. L., Kacham, P., and Woodruff, D. P. Near-optimal algorithms for linear algebra in the current matrix multiplication time. In *Proceedings of the 2022 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 3043–3068. SIAM, 2022.
- Chia, N., Lin, H., and Wang, C. Quantum-inspired sublinear classical algorithms for solving low-rank linear systems. *CoRR*, abs/1811.04852, 2018. URL http://arxiv.org/abs/1811.04852.
- Chia, N.-H., Gilyén, A., Li, T., Lin, H.-H., Tang, E., and Wang, C. Sampling-based sublinear low-rank matrix arithmetic framework for dequantizing quantum machine learning. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pp. 387–400, 2020.
- Chowdhury, A., Yang, J., and Drineas, P. An iterative, sketching-based framework for ridge regression. In *International Conference on Machine Learning*, pp. 989–998, 2018.
- Clarkson, K. L. Subgradient and sampling algorithms for  $\ell_1$  regression. In *Symposium on Discrete Algorithms: Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms*, volume 23, pp. 257–266, 2005.
- Clarkson, K. L. and Woodruff, D. P. Low rank approximation and regression in input sparsity time. In *STOC*, 2013. Full version at http://arxiv.org/abs/1207.6365. Final version J. ACM, Vol 63, 2017, http://doi.acm.org/10.1145/3019134.
- Clarkson, K. L. and Woodruff, D. P. Input sparsity and hardness for robust subspace approximation. In 2015 IEEE 56th Annual Symposium on Foundations of Computer Science, pp. 310–329. IEEE, 2015.
- Cohen, M. B. and Peng, R.  $L_p$  row sampling by Lewis weights. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pp. 183–192. ACM, 2015.
- Cohen, M. B., Elder, S., Musco, C., Musco, C., and Persu, M. Dimensionality reduction for k-means clustering and low rank approximation. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pp. 163–172, 2015.
- Cohen, M. B., Musco, C., and Musco, C. Input sparsity time low-rank approximation via ridge leverage score sampling. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1758–1777. SIAM, 2017.

- Cohen, M. B., Lee, Y. T., and Song, Z. Solving linear programs in the current matrix multiplication time. In *Proceedings of the 51st annual ACM SIGACT symposium on theory of computing*, pp. 938–942, 2019.
- Cong, I. and Duan, L. Quantum discriminant analysis for dimensionality reduction and classification. *New Journal* of *Physics*, 18(7):073011, 2016.
- Drineas, P., Mahoney, M. W., and Muthukrishnan, S. Subspace sampling and relative-error matrix approximation: Column-based methods. In *APPROX-RANDOM*, pp. 316–326, 2006.
- Drineas, P., Magdon-Ismail, M., Mahoney, M. W., and Woodruff, D. P. Fast approximation of matrix coherence and statistical leverage. *The Journal of Machine Learning Research*, 13(1):3475–3506, 2012.
- Dunjko, V. and Wittek, P. A non-review of quantum machine learning: trends and explorations. *Quantum Views*, 4:32, 2020.
- Frieze, A. M., Kannan, R., and Vempala, S. Fast Monte-Carlo algorithms for finding low-rank approximations. *J. ACM*, 51(6):1025–1041, 2004.
- Gilyén, A., Lloyd, S., and Tang, E. Quantum-inspired lowrank stochastic regression with logarithmic dependence on the dimension. *arXiv preprint arXiv:1811.04909*, 2018.
- Gilyén, A., Su, Y., Low, G. H., and Wiebe, N. Quantum singular value transformation and beyond: exponential improvements for quantum matrix arithmetics. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pp. 193–204, 2019.
- Gilyén, A., Song, Z., and Tang, E. An improved quantum-inspired algorithm for linear regression. *arXiv* preprint *arXiv*:2009.07268, 2020.
- Gilyén, A., Song, Z., and Tang, E. An improved quantum-inspired algorithm for linear regression. *CoRR*, abs/2009.07268, 2020.
- Harper, F. M. and Konstan, J. A. The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.*, 5(4):19:1–19:19, 2016. URL https://doi.org/10.1145/2827872.
- Harrow, A. W., Hassidim, A., and Lloyd, S. Quantum algorithm for linear systems of equations. *Physical review letters*, 103(15):150502, 2009.
- Indyk, P., Vakilian, A., Wagner, T., and Woodruff, D. Sample-optimal low-rank approximation of distance matrices. *arXiv preprint arXiv:1906.00339*, 2019.

- Jiang, H., Kathuria, T., Lee, Y. T., Padmanabhan, S., and Song, Z. A faster interior point method for semidefinite programming. arXiv preprint arXiv:2009.10217, 2020a.
- Jiang, S., Song, Z., Weinstein, O., and Zhang, H. Faster dynamic matrix inverse for faster LPs. *arXiv* preprint *arXiv*:2004.07470, 2020b.
- Kannan, R. and Vempala, S. Randomized algorithms in numerical linear algebra. *Acta Numerica*, 26:95, 2017.
- Kannan, R. and Vempala, S. S. Spectral algorithms. *Found. Trends Theor. Comput. Sci.*, 4(3-4):157–288, 2009.
- Kerenidis, I. and Prakash, A. Quantum recommendation systems. *arXiv preprint arXiv:1603.08675*, 2016.
- Laurent, B. and Massart, P. Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.*, 28 (5):1302–1338, 10 2000. doi: 10.1214/aos/1015957395. URL https://doi.org/10.1214/aos/1015957395.
- Lawrence, H., Li, J., Musco, C., and Musco, C. Low-rank toeplitz matrix estimation via random ultra-sparse rulers. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4796–4800. IEEE, 2020.
- Li, Y. and Woodruff, D. Input-sparsity low rank approximation in Schatten norm. *arXiv preprint arXiv:2004.12646*, 2020.
- Lloyd, S., Mohseni, M., and Rebentrost, P. Quantum principal component analysis. *Nature Physics*, 10(9):631–633, 2014.
- Lloyd, S., Garnerone, S., and Zanardi, P. Quantum algorithms for topological and geometric analysis of data. *Nature communications*, 7(1):1–7, 2016.
- Mahoney, M. W. Randomized algorithms for matrices and data. *Found. Trends Mach. Learn.*, 3(2):123–224, 2011.
- Musco, C. and Woodruff, D. P. Sublinear time low-rank approximation of positive semidefinite matrices. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 672–683. IEEE, 2017.
- Nelson, J. and Nguyên, H. L. OSNAP: Faster numerical linear algebra algorithms via sparser subspace embeddings. In 2013 ieee 54th annual symposium on foundations of computer science, pp. 117–126. IEEE, 2013.
- Rebentrost, P., Mohseni, M., and Lloyd, S. Quantum support vector machine for big data classification. *Phys. Rev. Lett.*, 113:130503, Sep 2014. doi: 10.1103/PhysRevLett.113. 130503. URL https://link.aps.org/doi/10.1103/PhysRevLett.113.130503.

- Rokhlin, V. and Tygert, M. A fast randomized algorithm for overdetermined linear least-squares regression. *Proceedings of the National Academy of Sciences*, 105(36): 13212–13217, 2008.
- Rudelson, M. and Vershynin, R. Sampling from large matrices: An approach through geometric functional analysis. *J. ACM*, 54(4), 2007.
- Sarlós, T. Improved approximation algorithms for large matrices via random projections. In *FOCS*, pp. 143–152, 2006.
- Shi, X. and Woodruff, D. P. Sublinear time numerical linear algebra for structured matrices. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4918–4925, 2019.
- Song, Z., Woodruff, D. P., and Zhong, P. Low rank approximation with entrywise  $\ell_1$ -norm error. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 688–701, 2017.
- Song, Z., Woodruff, D. P., and Zhong, P. Relative error tensor low rank approximation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 2772–2789. Society for Industrial and Applied Mathematics, 2019.
- Tang, E. A quantum-inspired classical algorithm for recommendation systems. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pp. 217–228. ACM, 2019.
- van Apeldoorn, J. and Gilyén, A. Improvements in Quantum SDP-Solving with Applications. In Baier, C., Chatzigiannakis, I., Flocchini, P., and Leonardi, S. (eds.), 46th International Colloquium on Automata, Languages, and Programming (ICALP 2019), volume 132 of Leibniz International Proceedings in Informatics (LIPIcs), pp. 99:1–99:15, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-109-2. doi: 10.4230/LIPIcs.ICALP.2019.99. URL http://drops.dagstuhl.de/opus/volltexte/2019/10675.
- Vishnoi, N. Cargese lecture notes. 2015. URL http://www.cs.yale.edu/homes/vishnoi/CargeseLectures.pdf.
- Wang, R. and Woodruff, D. P. Tight bounds for  $\ell_p$  oblivious subspace embeddings. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1825–1843. SIAM, 2019.
- Woodruff, D. P. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.

Zhao, Z., Fitzsimons, J. K., and Fitzsimons, J. F. Quantum-assisted gaussian process regression. *Phys. Rev. A*, 99:052331, May 2019. doi: 10.1103/PhysRevA.99.052331. URL https://link.aps.org/doi/10.1103/PhysRevA.99.052331.

#### A. Preliminaries

In, this section, we provide proofs for our theorems in Section 3. We first provide the description of the leverage score sampling data structure (Algorithm 5). We note that the advantage of the current version compared to the standard leverage score sampling (see, e.g., Section 2.4 in the survey in (Woodruff, 2014)) is that it saves an  $O(\log n)$  factor. We need the following data structure.

**Definition A.1** (Sampling Data Structure). Given a matrix  $A \in \mathbb{R}^{n \times d}$ , a column selection matrix  $\Lambda$  such that  $k = \text{rank}(A\Lambda) = \text{rank}(A)$ , and  $\lambda_s \geq 1$ , the data structure SAMP $(A, \Lambda, \lambda_s)$  consists of the following:

- SA, where  $S \in \mathbb{R}^{m_S \times n}$  is a sketching matrix as in Lemma 3.1, with  $m_S = O(k \log(k)/\varepsilon_0^2)$ , chosen to be an  $\varepsilon_0$ -embedding with failure probability 1/100, for fixed  $\varepsilon_0$ ;
- C, where  $[Q, C] \leftarrow QR(SA\Lambda)$ , the QR decomposition of  $SA\Lambda$ , i.e.,  $SA\Lambda = QC$ , Q has orthonormal columns and C is triangular and invertible (since  $A\Lambda$  has full column rank);
- $C_0$ , where  $[Q_0, C_0] \leftarrow \mathsf{QR}(SA)$ ;
- The data structure of Lemma 4.1, built to enable sampling  $i \in [n]$  with probability  $p_i \leftarrow \|Z_{i,*}\|^2 / \|Z\|_F^2$  in  $O(\log n)$  time, where  $Z \leftarrow A\Lambda(C^{-1}G)$ , with  $G \in \mathbb{R}^{k \times m_G}$  having independent  $\mathcal{N}(0, 1/m_G)$  entries, and  $m_G = \Theta(\lambda_s)$ .

#### **Algorithm 4** MATVECSAMPLER(A, SAMPLER, W, v, $\nu$ )

Input:  $A \in \mathbb{R}^{n \times d}$ , data structure SAMPLER (Def. A.1),  $W \in \mathbb{R}^{d \times m_W}$ , desired number of samples v, normalizer v, where  $v = \frac{1}{6kn^{1/\lambda_s}}$  by default if unspecified

**Output:**  $L \in \mathbb{R}^{v \times d}$ , encoding v draws from  $i \in [n]$  chosen with approx. probability  $q_i \stackrel{def}{=} \|A_{i,*}W\|^2 / \|AW\|_F^2$ 

```
1: N \leftarrow \|C_0 W\|_F^2, where C_0 is from Sampler

2: if N = 0:

3: return Uniform(v, n) {Alternatively, raise an exception here}

4: L \leftarrow 0_{v \times n}, z \leftarrow 0

5: while z < v:

6: Choose i \in [n] with probability p_i using Sampler

7: \tilde{q}_i \leftarrow \|A_{i,*}W\|^2/N

8: With probability \nu^{\tilde{q}_i}_{p_i}, accept i: set L_{z,i} = 1/\sqrt{v\tilde{q}_i}; z \leftarrow z + 1

9: return L
```

**Lemma A.2** (Sampling Data structure). The data structure SAMP $(A, \Lambda, \lambda_s)$ , from Definition A.1, can be constructed in  $O(\lambda_s(\text{nnz}(A) + k^2) + d^{\omega})$  time.

*Proof.* The time needed to compute SA is  $O(\operatorname{nnz}(A) + k^{\omega}\operatorname{poly}(\log\log(k)) + k^{2+o(1)}\varepsilon_0^{-2})$ . Computing the QR factorization of SA takes  $O(d^{\omega})$  time, by first computing  $(SA)^{\top}(SA)$  for the  $m_S \times d$  matrix SA, and then its Cholesky composition, using "fast matrix" methods for both, and using  $m_S \leq d$ . This dominates the time for the similar factorization of  $SA\Lambda$ .

The Z matrix can be computed in  $O(\lambda_s(\mathtt{nnz}(A) + k^2))$  time, by appropriate order of multiplication, and this dominates the time needed for building the data structure of Lemma 4.1. Adding these terms, and using  $m \leq \mathtt{nnz}(A)$ , the result follows.

**Lemma A.3** (MatVecSampler Analysis). Given constant  $c_0 > 1$  and small enough constant  $\varepsilon_0 > 0$ , and Sampler for A, there is an event  $\mathcal E$  holding with failure probability at most  $1/k^{c_0}$ , so that if  $\mathcal E$  holds, then when called with  $\nu \leftarrow \frac{1}{6kn^{1/\lambda_s}}$ , the probability is  $(1\pm \varepsilon_0)q_i$  that the accepted index in Step 4 of MatVecSampler is  $i\in [n]$ , where  $q_i\stackrel{def}{=}\|A_{i,*}W\|^2/\|AW\|_F^2$ . The time taken is  $O(m_W d(d+vkn^{1/\lambda_s}))$ , where  $k=\mathrm{rank}(A)$ .

*Proof.* We need to verify that the quantity in question is a probability, that is, that  $\nu \frac{\tilde{q}_i}{p_i} \in (0,1)$ , when  $\nu = \frac{1}{6kn^{1/\lambda_s}}$ .

From Lemma 3.1, if  $m_S = O(\varepsilon_0^{-2}k)$  for  $\varepsilon_0 > 0$ , then with failure probability  $1/k^{c_0+1}$ , S will be a subspace  $\varepsilon_0$ -embedding for  $\operatorname{im}(A)$ , that is, for  $A\Lambda$  and for A, using  $\operatorname{rank}(A) = k$ . The event  $\mathcal E$  includes the condition that S is

indeed an  $\varepsilon_0$ -embedding. If this holds for S, then from standard arguments,  $A\Lambda C^{-1}$  has singular values all in  $1 \pm \varepsilon_0$ , and  $\|A_{i,*}\Lambda C^{-1}\|^2 = (1 \pm O(\varepsilon_0))\tau_i$ , where again  $\tau_i$  is the i'th leverage score.

(We have  $||A\Lambda C^{-1}x|| = (1 \pm \varepsilon_0)||SA\Lambda C^{-1}x|| = (1 \pm \varepsilon_0)||x||$ , for all x, since SA = QC.) This implies that for Z, G in the construction of SAMP $(A, \Lambda, \lambda_s)$ ,

$$||Z||_F^2 = ||A\Lambda C^{-1}G||_F^2 \le (1+\varepsilon_0)||G||_F^2.$$

Since  $m_G \|G\|_F^2$  is  $\chi^2$  with  $km_G$  degrees of freedom, with failure probability at most  $\exp(-\sqrt{k\lambda_s}/2)$  (using  $m_G = \Theta(\lambda_s)$ , it is at most  $3km_G$  ((Laurent & Massart, 2000), Lemma 1), so  $\|G\|_F^2 \leq 3k$  with that probability. Our event  $\mathcal E$  also includes the condition that this bound holds. Thus under this condition,  $\|Z\|_F^2 \leq 3(1+\varepsilon_0)k$ .

From Lemma A.8 and the above characterization of  $\tau_i$ , for the Z of SAMP $(A, \Lambda, \lambda_s)$ ,  $\|Z_{i,*}\|^2 = \|A_{i,*}\Lambda C^{-1}G\|^2 \ge (1 - O(\varepsilon_0))\tau_i/n^{1/\lambda_s}$ .

Putting these together, we have

$$p_{i} = \frac{\|Z_{i,*}\|_{F}^{2}}{\|Z\|^{2}} \ge (1 - O(\varepsilon_{0})) \frac{\tau_{i}/n^{1/\lambda_{s}}}{3k}.$$
(1)

Using the  $\varepsilon_0$ -embedding property of S,

$$||C_0W||_F^2 = ||Q_0C_0W||_F^2 = ||SAW||_F^2 = (1 \pm 2\varepsilon_0)||AW||_F^2,$$
(2)

and so, letting  $A = UC_1$  for U with orthonormal columns, we have, for small enough  $\varepsilon_0$ ,

$$(1 - 2\varepsilon_0)\tilde{q}_i \le \frac{\|A_{i,*}W\|^2}{\|AW\|_F^2} = \frac{\|U_{i,*}C_1W\|^2}{\|UC_1W\|_F^2} = \frac{\|U_{i,*}C_1W\|^2}{\|C_1W\|_F^2} \le \frac{\|U_{i,*}\|^2\|C_1W\|^2}{\|C_1W\|_F^2} \le \tau_i.$$

Putting this bound with (1) we have

$$\frac{\tilde{q}_i}{p_i} \le \frac{\tau_i/(1-2\varepsilon_0)}{(1-O(\varepsilon_0))\tau_i/n^{1/\lambda_s}3(1+\varepsilon_0)k} \le 3(1+O(\varepsilon_0))kn^{1/\lambda_s}.$$

so that  $\nu \frac{\tilde{q}_i}{p_i} = \frac{1}{6kn^{1/\lambda_s}} \frac{\tilde{q}_i}{p_i} \le (1 + O(\varepsilon_0))/2 \le 1$  for small enough  $\varepsilon_0$ . Using (2) we have  $\tilde{q}_i = (1 \pm 2\varepsilon_0)q_i$ . Thus the correctness condition of the lemma follows, for small enough  $\varepsilon_0$ .

Turning to time: the time to compute  $C_0W$  is  $O(d^2m_W)$ . Each iteration takes  $O(\log n + dm_W)$ , for choosing i and computing  $\tilde{q}_i$ , and these steps dominate the time. As usual for rejection sampling, the expected number of iterations is  $O(vkn^{1/\lambda_s})$ . Adding these expressions yields the expected time bound, folding a factor of  $\log n$  in by adjusting  $\lambda_s$  slightly.

## Algorithm 5 LEVSAMPLE $(A, \mu_s, f())$

Input:  $A \in \mathbb{R}^{n \times d}$ ,  $\mu_s \ge 1$  specifying runtime tradeoff, function  $f(\cdot) \to \mathbb{Z}_+$  returns a target sample size (may be just a constant) Output: Leverage score sketching matrix L, column selector  $\Lambda$ 

- 1: Run an algorithm to compute  $k = \mathtt{rank}(A)$  and obtain  $\Lambda \in \mathbb{R}^{d \times k}$ , a subset of k lin. indep. columns of A {for example Theorem 1.5 in (Chepurko et al., 2022)}
- 2: Construct SAMPLER  $\leftarrow$  SAMP $(A\Lambda, I, \lambda_s)$ , use C from it; {Definition A.1}
- 3:  $W \leftarrow C^{-1}G'$ , where  $G' \in \mathbb{R}^{k \times m_{G'}}$  with ind.  $\mathcal{N}(0, 1/m_{G'})$  entries  $\{m_{G'} = \Theta(\log n)\}$
- 4:  $L \leftarrow \text{MATVECSAMPLER}(A\Lambda, \text{SAMPLER}, W, f(k, C), \nu = 1/6n^{1/\lambda_s})$ {Algorithm 4, sample size f(k, C), normalizer  $\nu$ }
- 5: return  $L, \Lambda$

**Theorem A.4** (Leverage Score Data Structure, Theorem 3.2 restated). Let k = rank(A), and choose  $\mu_s \ge 1$ . Algorithm 5 (LEVSAMPLE $(A, \mu_s, f(\cdot))$ ) uses space  $O(n + k^{\omega} \log \log(nd))$ , not counting the space to store A, and runs in time

$$O(\mu_s \operatorname{nnz}(A) + k^{\omega} \operatorname{poly}(\log \log(k)) + k^{2+o(1)} + v k n^{1/\mu_s}),$$

where v is the sample size. For  $v = \varepsilon^{-2}k \log k$ , this bound can be expressed as  $O(\mu_e \operatorname{nnz}(A) + k^\omega \operatorname{poly}(\log \log(k)) + \varepsilon^{-2-1/\mu_e}k^2)$  time, for  $\mu_e \geq 1$ .

Note: we can get a slightly smaller running time by including more rounds of rejection sampling: the first round of sampling needs an estimate with failure probability totaled over for all n rows, while another round would only need such a bound for  $vn^{1/\lambda_s}$  rows; this would make the bound  $v^{1+1/\lambda_s}kn^{1/\lambda_s^2}$ , which would be smaller when  $v \ll n$ . However, in the latter case, the term  $vkn^{1/\lambda_s}$  is dominated by the other terms anyway, for relevant values of the parameters. For example if  $v \leq n$  and  $vk \leq \text{nnz}(A)$  does not hold, then sampling is not likely to be helpful. Iterating  $\log \log n$  times, a bound with leading term  $O(\text{nnz}(A)(\log \log n + \log v))$  is possible, but does not seem interesting.

*Proof.* Step A, building SAMP $(A\Lambda, \lambda_s)$ , take  $O(\lambda_s(\text{nnz}(A) + k^2) + k^\omega)$  time, with d in Lemma A.2 equal to k here.

From Lemma A.3, the running time of MATVECSAMPLER is  $O(k^2 \log n + vk^2(\log n)n^{1/\lambda_s})$ , mapping d of the lemma to k,  $m_W$  to  $m_{G'} = O(\log n)$ . However, since the normalizer  $\nu$  is a factor of k smaller than assumed in Lemma A.3, the runtime in sampling is better by that factor. Also, we subsume the second  $\log n$  factor by adjusting  $\lambda_s$ .

The cost of computing  $C^{-1}G'$  is  $O(k^2 \log n)$ ; we have a runtime of

$$\begin{split} O(\mathtt{nnz}(A) + k^\omega \mathsf{poly}(\log\log(k)) + k^{2+o(1)}) + O(\lambda_s(\mathtt{nnz}(A) + k^2) + k^\omega) + O(k^2 \log n + vkn^{1/\lambda_s}) \\ &= O(\lambda_s \, \mathtt{nnz}(A) + k^\omega \mathsf{poly}(\log\log(k)) + k^{2+o(1)} + vkn^{1/\lambda_s}), \end{split}$$

as claimed.

Finally, suppose  $v = \varepsilon^{-2}k \log k$ , as suffices for an  $\varepsilon$ -embedding. If  $vkn^{1/\lambda_s} \leq \operatorname{nnz}(A) + k^{\omega}$ , then the bound follows. Suppose not. If  $n \geq k^{\omega}$ , then

$$\varepsilon^{-2} \ge n^{1-1/\lambda_s}/k^2 \log(k) \ge n^{1-1/\lambda_s-2/\omega}/\log(n)$$

and so  $\varepsilon^{-2} \ge n^{\gamma}$ , for constant  $\gamma > 0$ , implying  $\varepsilon^{-2/\lambda_s \gamma'} \ge n^{1/\lambda_s} \log n$ , for constant  $\gamma' < \gamma$ . When  $k^{\omega} \ge n$ ,

$$\varepsilon^{-2} \ge k^{\omega - 2 - 1/\omega \lambda_s} / \log(k) \ge k^{\gamma},$$

for a constant  $\gamma > 0$ , so that  $\varepsilon^{-\omega/\lambda_s \gamma'} \ge n^{1/\lambda_s} \log k$ , for a constant  $\gamma' < \gamma$ . Using  $\lambda_e$ , a constant multiple of  $\lambda_s$ , to account for constants, the result follows.

We can also use *length-squared sampling* to obtain subspace embeddings. In Section 4 we have given a data structure and algorithm that implements length-squared sampling. To analyze length-squared sampling in the context of ridge regression, we show the following structural observations about ridge regression.

**Lemma A.5** (Block SVD). Let A be such that  $k = \operatorname{rank}(A)$ , and suppose A has thin SVD  $A = U\Sigma V^{\top}$ , implying  $\Sigma \in \mathbb{R}^{k \times k}$ . For  $\lambda > 0$ , let  $A_{(\lambda)} = \begin{bmatrix} A \\ \sqrt{\lambda}VV^{\top} \end{bmatrix}$ . For  $b \in \mathbb{R}^n$ , let  $\hat{b} = \begin{bmatrix} b \\ 0d \end{bmatrix}$ . Then for all  $x \in \operatorname{im}(V)$ , the ridge regression loss

$$||Ax - b||^2 + \lambda ||x||^2 = ||A_{(\lambda)}x - \hat{b}||^2,$$

and ridge regression optimum

$$x^* = \operatorname{argmin}_{x \in \mathbb{R}^d} \|Ax - b\|^2 + \lambda \|x\|^2 = \operatorname{argmin}_{x \in \mathbb{R}^d} \|A_{(\lambda)}x - \hat{b}\|^2$$

The matrix  $A_{(\lambda)}$  has SVD  $A_{(\lambda)} = \begin{bmatrix} U\Sigma D \\ \sqrt{\lambda}VD \end{bmatrix} D^{-1}V^{\top}$ , where  $D = (\Sigma^2 + \lambda I_k)^{-1/2}$ , and  $\|A_{(\lambda)}^+\|^2 = 1/(\lambda + 1/\|A^+\|^2)$ . We have  $\|A_{i,*}\|^2 \|A_{(\lambda)}^+\|^2 \ge \|U_{i,*}\Sigma D\|^2$  for  $i \in [n]$ .

*Proof.* Since  $x \in \text{im}(V)$  has x = Vz for some  $z \in \mathbb{R}^k$ , and since  $V^\top V = I_k$ , it follows that  $VV^\top x = Vz = x$ , and so

$$||A_{(\lambda)}x - \hat{b}||^2 = ||Ax - b||^2 + ||\sqrt{\lambda}VV^{\top}x - 0||^2 = ||Ax - b||^2 + \lambda||x||^2,$$

as claimed.

The SVD of  $A_{(\lambda)}$  is  $A_{(\lambda)} = \begin{bmatrix} U\Sigma D \\ \sqrt{\lambda}VD \end{bmatrix} D^{-1}V^{\top}$ , where D is defined as in the lemma statement, since the equality holds, and both  $\begin{bmatrix} U\Sigma D \\ \sqrt{\lambda}VD \end{bmatrix}$  and V have orthonormal columns, and  $D^{-1}$  has non-increasing nonnegative entries. Therefore  $A_{(\lambda)}^+ = VD \begin{bmatrix} U\Sigma D \\ \sqrt{\lambda}VD \end{bmatrix}^{\top}$ . We have

$$A_{(\lambda)}^{+}\hat{b} = VD^{2}\Sigma U^{\top}b = V\Sigma D^{2}U^{\top}b, \tag{3}$$

using that  $\Sigma$  and D are diagonal matrices.

By the well-known expression  $x^* = A^{\top} (AA^{\top} + \lambda I_n)^{-1} b$ , and using the *not*-thin SVD  $A = \hat{U} \hat{\Sigma} \hat{V}^{\top}$ , with  $\hat{\Sigma} \in \mathbb{R}^{n \times d}$  and  $\hat{U}$  and  $\hat{V}$  orthogonal matrices,

$$x^* = \hat{V}\hat{\Sigma}\hat{U}^{\top}(\hat{U}\hat{\Sigma}\hat{\Sigma}^{\top}\hat{U}^{\top} + \lambda\hat{U}\hat{U}^{\top})^{-1}b$$

$$= \hat{V}\hat{\Sigma}\hat{U}^{\top}\hat{U}(\hat{\Sigma}\hat{\Sigma}^{\top} + \lambda I_n)^{-1}\hat{U}^{\top}b$$

$$= \hat{V}\hat{\Sigma}(\hat{\Sigma}\hat{\Sigma}^{\top} + \lambda I_n)^{-1}\hat{U}^{\top}b$$

$$= V\Sigma(\Sigma^2 + \lambda I_k)^{-1}U^{\top}b$$

$$= V\Sigma D^2 U^{\top}b,$$

where the next-to-last step uses that  $\hat{\Sigma}$  is zero except for the top k diagonal entries of  $\hat{\Sigma}$ . Comparing (3) and (4), we have  $A_{(\lambda)}^+\hat{b}=x^*$ . Using the expression for  $A_{(\lambda)}^+$ ,  $\|A_{(\lambda)}^+\|^2=D_{1,1}^2=1/(\lambda+1/\|A^+\|^2)$ . Finally, since  $(A_{(\lambda)})_{i,*}=A_{i,*}$  for  $i\in[n]$ , and letting  $\hat{U}=\begin{bmatrix}U\Sigma D\\\sqrt{\lambda}VD\end{bmatrix}$ ,

$$||A_{i,*}||^2 ||A_{(\lambda)}^+||^2 = ||(A_{(\lambda)})_{i,*}||^2 ||A_{(\lambda)}^+||^2$$

$$\geq ||(A_{(\lambda)})_{i,*}A_{(\lambda)}^+||^2$$

$$= ||\hat{U}_{i,*}D^{-1}V^\top V D \hat{U}^\top||^2$$

$$= ||\hat{U}_{i,*}||^2$$

$$= ||U_{i,*}\Sigma D||^2,$$

as claimed.

**Definition A.6** (Ridge Leverage-score Sample, Statistical Dimension). Let  $A, \lambda, A_{(\lambda)}$ , and D be as in Lemma A.5. Call  $\mathcal{S} \subset [n]$  a *ridge leverage-score sample* of A if each  $i \in \mathcal{S}$  is chosen independently with probability at least  $\|U_{i,*}\Sigma D\|^2/\operatorname{sd}_{\lambda}(A)$ , where the statistical dimension  $\operatorname{sd}_{\lambda}(A) = \|U\Sigma D\|_F^2 = \sum_{i \in [d]} \sigma_i^2/(\lambda + \sigma_i^2)$ , recalling that  $U\Sigma D$  comprises the top n rows of the left singular matrix of  $A_{(\lambda)}$ .

**Lemma A.7** (Length-squared sketch, Lemma 3.5 restated). Given a matrix  $A \in \mathbb{R}^{n \times d}$  and a sample size parameter  $v \in [n]$ , let  $m = O\left(v\|A_{(\lambda)}^+\|^2\|A\|_F^2/\operatorname{sd}_{\lambda}(A)\right)$ . Then, with probability at least 99/100, the set of m length-squared samples contains a ridge leverage-score sample of A of size v.

Note that when  $\lambda = 0$ ,  $||A_{(\lambda)}^+|| = ||A^+||$ ,  $sd_0(A) = rank(A)$ , and the ridge leverage-score samples are leverage-score samples.

*Proof.* We will show that  $\hat{L}$  contains within it a leverage-score sketching matrix; since oversampling does no harm, this implies the result using the above lemma.

Using the thin SVD  $A = U\Sigma V^{\top}$ , and  $A^+ = V\Sigma^+ U^{\top}$ , we have

$$||A_{i,*}|| ||A^+|| \ge ||A_{i,*}A^+|| = ||U_{i,*}\Sigma V^\top V \Sigma^+ U^\top|| = ||U_{i,*}||,$$

The expected number of times index  $i \in [n]$  is chosen among  $m_{\hat{L}}$  length-squared samples,  $p_i m_{\hat{L}}$ , is within a constant factor of  $\frac{\|A_{i,*}\|^2}{\|A\|_F^2} v \|A_{(\lambda)}^+\|^2 \|A\|_F^2 / \operatorname{sd}_{\lambda}(A) \geq \frac{\|U_{i,*}\|^2}{\operatorname{sd}_{\lambda}(A)} v$ , using Lemma A.5, an expectation at least as large as for a ridge leverage-score sample of size v.

Finally, recall the Johnson-Lindenstraus Lemma, for sketching with a dense Gaussian matrix.

**Lemma A.8** (Johnson-Lindenstraus Lemma). For given  $\varepsilon > 0$ , if  $P \subset \mathbb{R}^c$  is a set of  $m \ge c$  vectors, and  $G \in \mathbb{R}^{m \times c}$  has entries that are independent Gaussians with mean zero and variance 1/m, then there is  $m = O(\varepsilon^{-2} \log(m/\delta))$  such that with failure probability  $\delta$ ,  $||Gx|| = (1 \pm \varepsilon)||x||$  for all  $x \in P$ . Moreover, there is  $m_G = O(\mu)$  so that  $||Gx|| \ge ||x||/n^{1/\mu}$ , with failure probability at most  $1/n^2$ .

**Definition A.9** (Projection-Cost Preserving Sketch). Given a matrix  $A \in \mathbb{R}^{n \times d}$ ,  $\varepsilon > 0$  and an integer  $k \in [d]$ , a sketch  $SA \in \mathbb{R}^{s \times d}$  is an  $(\varepsilon, k)$ -column projection-cost preserving sketch of A if for all rank-k projection matrices P,  $(1-\varepsilon)\|A(I-P)\|_F^2 \leq \|SA(I-P)\|_F^2 \leq (1+\varepsilon)\|A(I-P)\|_F^2$ .

There are several constructions of projection-cost preserving sketches known in the literature, starting with the work of Cohen et al., 2015; 2017). For our purposes, it suffices to use Theorem 1 from (Cohen et al., 2017).

We can use the following lemma to translate from prediction error to solution error for regression problems.

**Lemma A.10.** Let  $\gamma_{A,b} = \frac{\|b\|}{\|AA^+b\|}$ . Recall that  $\kappa(A) = \|A\| \|A^+\|$ . Suppose  $\tilde{x} \in \text{im}(A^\top)$ , and for some  $\varepsilon_p \in (0,1)$ ,  $\|A\tilde{x} - b\|^2 \le (1 + \varepsilon_p) \|\xi^*\|^2$  holds, where  $\xi^* = Ax^* - b$ . Then

$$\|\tilde{x} - x^*\| \le 2\sqrt{\varepsilon_p}\|A^+\|\|\xi^*\| \le 2\sqrt{\varepsilon_p}\sqrt{\gamma_{A,b}^2 - 1}\kappa(A)\|x^*\|.$$
 (4)

This extends to multiple response regression using  $\gamma_{A,B}^2 \stackrel{\text{def}}{=} \frac{\|B\|_F^2}{\|AA+B\|_F^2}$ , by applying column by column to B, and extends to ridge regression, that is,  $A_{(\lambda)}$  with  $\hat{B} = \begin{bmatrix} B \\ 0_{d \times d'} \end{bmatrix}$ , as well.

Note that  $x \in \operatorname{im}(A^{\top}) = \operatorname{im}(V)$  is no loss of generality, because the projection  $VV^{\top}x$  of x onto  $\operatorname{im}(A^{\top})$  has  $AVV^{\top}x = Ax$  and  $\|VV^{\top}x\| \leq \|x\|$ . So  $\operatorname{argmin}_x \|Ax - b\|^2 + \lambda \|x\|$  must be in  $\operatorname{im}(A^{\top})$  for  $\lambda$  arbitrarily close to zero, and  $A^+b \in \operatorname{im}(A^{\top})$ .

For the ridge problem  $\min_x \|A_{(\lambda)}x - \hat{b}\|$ , we have  $\|\xi^*\|^2 = \|A_{(\lambda)}x^* - \hat{b}\| = \|Ax^* - b\|^2 + \lambda \|x^*\|^2$ , and recalling from Lemma A.5 that, when A has SVD  $A = U\Sigma V^\top$ ,  $A_{(\lambda)}$  has singular value matrix  $D^{-1}$ , where  $D = (\Sigma^2 + \lambda I_k)^{-1/2}$ , so that  $\kappa(A_{(\lambda)})^2 = (\lambda + \sigma_1^2)/(\lambda + \sigma_k^2)$ , where A has singular values  $\sigma_1, \ldots, \sigma_k$ , with k = rank(A).

*Proof.* Since  $x^* = A^+b = A^\top (AA^\top)^+b \in \operatorname{im} A^\top$ , we have  $\tilde{x} - x^* = A^\top z \in \operatorname{im} A^\top$ , for some z. Since  $A^+AA^\top = A^\top$ , we have  $\tilde{x} - x^* = A^\top z = A^+AA^\top z = A^+A(\tilde{x} - x^*)$ . From the normal equations for regression and the Pythagorean theorem,

$$||A(\tilde{x} - x^*)||^2 = ||A\tilde{x} - b||^2 - ||Ax^* - b||^2 \le 4\varepsilon_p ||\xi^*||^2,$$

using  $||A\tilde{x} - b|| \le (1 + \varepsilon_p)||\xi^*||$  and  $\varepsilon_p < 1$ . Therefore, using also submultiplicativity of the spectral norm,

$$\|\tilde{x} - x^*\|^2 = \|A^+ A(\tilde{x} - x^*)\|^2$$

$$\leq \|A^+\|^2 \|A(\tilde{x} - x^*)\|^2$$

$$\leq \|A^+\|^2 4\varepsilon_p \|\xi^*\|^2,$$
(5)

and the first inequality of (4) follows. For the second, we bound

$$\frac{\left\|\xi*\right\|^{2}}{\left\|x^{*}\right\|^{2}} = \frac{\left\|b\right\|^{2} - \left\|AA^{+}b\right\|^{2}}{\left\|A^{+}b\right\|^{2}} = \frac{(\gamma_{A,b}^{2} - 1){\left\|AA^{+}b\right\|^{2}}}{\left\|A^{+}b\right\|^{2}} \le (\gamma_{A,b}^{2} - 1){\left\|A\right\|^{2}}$$

so from (5), we have

$$\|\tilde{x} - x^*\|^2 \le \|A^+\|^2 4\varepsilon_p \|\xi^*\|^2 \le \|A^+\|^2 4\varepsilon_p \|x^*\|^2 (\gamma_{A,b}^2 - 1) \|A\|^2$$

and the second inequality of (4) follows, using the definition of  $\kappa(A)$ .

### **B. Dynamic Data Structures for Ridge Regression**

The data structure of Lemma 4.1 is simply a complete binary tree with  $\ell$  leaves, each leaf with weight  $u_i$ , and each internal node with weight equal to the sum of the weights of its children. Sampling is done by walking down from the root, choosing left or right children with probability proportional to its weight. Insertion and deletion are done by inserting or deleting the leaf z that preserves the complete binary tree property, and updating the weights of its ancestors; in the case of deletion, first the leaf weight to be deleted is swapped with that of z, updating weights of ancestors. We also refer the reader to a more detailed description in (Tang, 2019; Gilyén et al., 2020).

Proof of Lemma 4.3. Use Lemma 4.1 for the first part. For the second, with a matrix S, construct the data structure of Lemma 4.1 for the row lengths of SA, in  $O(m\log(nd))$  time. To sample, pick  $i^* \in [m]$  with probability  $\|(SA)_{i,*}\|^2/\|SA\|_F^2$ , using the newly constructed data structure. Then pick  $j \in [d]$  with probability  $(SA)_{i^*j}/\|(SA)_{i^*,*}\|^2$ . Adding the probabilities across the choices of  $i^*$ , the probability of choosing index j is  $\|(SA)_{*,j}\|^2/\|SA\|_F^2$ , as claimed.

Once a column is chosen, the time to determine the corresponding column length  $||(SA)_{*,j}||$  is  $O(m \log(nd))$ , finding each  $(SA)_{ij}$  for  $i \in [m]$  in  $O(\log(nd))$  time.

We first re-state our theorem for dynamic ridge regression, before giving its proof.

**Theorem B.1** (Theorem 1.4 restated). Given matrices  $A \in \mathbb{R}^{n \times d}$ ,  $B \in \mathbb{R}^{n \times d'}$  and  $\lambda > 0$ , let  $X^* = \underset{K}{\operatorname{argmin}}_X \|AX - B\|_F^2 + \lambda \|X\|_F^2$ . Let  $\psi_{\lambda} = \|A\|_F^2/(\lambda + \sigma_k^2)$ ,  $\kappa_{\lambda} = (\lambda + \sigma_1(A)^2)/(\lambda + \sigma_k(A)^2)$  and  $\hat{\kappa} = \sqrt{(\lambda + \hat{\sigma}_1^2)/(\lambda + \hat{\sigma}_k^2)}$ , where  $\hat{\sigma}_1$  and  $\hat{\sigma}_k$  are over and under estimates of  $\sigma_1$  and  $\sigma_k$  respectively. Then, there exists a data structure that maintains  $Y \in \mathbb{R}^{d \times d'}$  such that with probability at least 99/100,

$$\|Y - X^*\|_F \le (\varepsilon + 2\gamma\varepsilon) \|X^*\|_F + \frac{\varepsilon}{\sqrt{\lambda}} \|U_{k,\perp}B\|_F,$$

where  $\gamma^2 = \frac{\|B\|_F^2}{\|AA^+X^*\|_F^2}$ . Further, an entry  $Y_{ij}$  for given i,j can be computed in  $O(m_S \log(nd)) = O(\varepsilon^{-2}\hat{\kappa}^2\psi_\lambda(\log(nd))^2)$  time. The time taken to compute Y is  $\tilde{O}(d'\varepsilon^{-4}\hat{\kappa}^2\psi_\lambda^2\kappa_\lambda\log(d))$ .

Proof. Let

$$X_1 = \operatorname{argmin}_{X \in \mathbb{R}^{d \times d'}} \|SAX - SB\|_F^2 + \lambda \|X\|_F^2.$$

We first show that

$$||X_1 - X^*||_F \le \varepsilon \gamma_{A_{(\lambda)}, \hat{B}} ||X^*||_F, \tag{6}$$

which follows from Lemma A.10, applied to  $A_{(\lambda)}$  and  $\hat{B}$ , after showing that, for  $\varepsilon_p = \varepsilon^2/\hat{\kappa}^2$ ,  $X_1$  satisfies

$$||AX_1 - B||_F^2 + \lambda ||X_1||_F^2 \le (1 + \varepsilon_p/4)\Delta_*, \text{ where } \Delta_* = ||AX^* - B||_F^2 + \lambda ||X^*||_F^2,$$
(7)

which in turn follows from Lemma 17 of (Avron et al., 2017). That lemma considers a matrix  $U_1$ , comprising the first n rows of the left singular matrix of  $\hat{A}_{(\lambda)} = \begin{bmatrix} A \\ \sqrt{\lambda}I_d \end{bmatrix}$ , noting that the ridge objective can be expressed as  $\min_X \|\hat{A}_{(\lambda)}X - [^B_0]\|_F^2$ . The matrix  $U_1 = U\Sigma D$  in our terminology, as in Lemma A.5, so the observations of that lemma apply.

Lemma 17 of (Avron et al., 2017) requires that S satisfies

$$||U_1^{\top} S^{\top} S U_1 - U_1^{\top} U_1|| \le 1/4, \tag{8}$$

and

$$||U_1^\top S^\top S(B - AX^*) - U_1^\top (B - AX^*)||_F \le \sqrt{\varepsilon_p \Delta_*}.$$
(9)

We have  $\|A_{(\lambda)}^+\|^2 = 1/(\lambda + 1/\|A^+\|^2) \le Z_\lambda$ . With the given call to LenSqSample to construct S, the number of rows sampled is  $m_S = O(\varepsilon_p^{-1} Z_\lambda^2 \|A\|_F^2 \log(d))$ , so the expected number of times that row i of A is sampled is, up to a factor of  $O(\log d)$ ,

$$\varepsilon_p^{-1} Z_{\lambda}^2 \|A\|_F^2 \frac{\|A_{i,*}\|^2}{\|A\|_F^2} = \varepsilon_p^{-1} Z_{\lambda}^2 \|A_{i,*}\|^2 \ge \varepsilon_p^{-1} \|(U_1)_{i,*}\|^2 = \varepsilon_p^{-1} \|U_1\|_F^2 \frac{\|(U_1)_{i,*}\|^2}{\|U_1\|_F^2},$$

and so row i is sampled at least the expected number of times it would be sampled under  $\varepsilon_p^{-1} ||U_1||_F^2 \log d$  rounds of length-squared sampling of  $U_1$ . As shown by Rudelson and Vershynin ((Rudelson & Vershynin, 2007), see also (Kannan & Vempala, 2017), Theorem 4.4), this suffices to have, with high probability, a bound on the normed expression in (8) of

$$\frac{\|U_1\|\|U_1\|_F}{\sqrt{\varepsilon_p^{-1}\|U_1\|_F^2}} = \sqrt{\varepsilon_p}\|U_1\| \le \sqrt{\varepsilon_p},$$

so by adjusting constant factors in sample size, (8) holds, for small enough  $\varepsilon_p$ .

To show that (9) holds, we use the discussion of the basic matrix multiplication algorithm discussed in (Kannan & Vempala, 2017), Section 2.1, which implies that

$$E[\|U_1^{\top} S^{\top} S(B - AX^*) - U_1^{\top} (B - AX^*)\|_F^2] \le \frac{\|U_1\|_F^2 \|B - AX^*\|_F^2}{s}$$

where s is the number of length-squared samples of  $U_1$ . Here  $s = \varepsilon_p^{-1} ||U_1||_F^2 \log d$ , so (9) follows with constant probability by Chebyshev's inequality, noting that  $||B - AX^*||_F \le \sqrt{\Delta_*}$ .

Thus (8) and (9) hold, so that by Lemma 17 of (Avron et al., 2017), (7) holds. We now apply Lemma A.10, which with (7) and  $\varepsilon_p = \varepsilon^2/\hat{\kappa}^2$ , implies (6).

Next we show that the (implicit) returned solution is close to the solution of (7), that is,

$$\|A^{\top}S^{\top}\tilde{X} - X_1\|_F^2 \le \varepsilon \|X_1\|_F^2. \tag{10}$$

This is implied by Theorem 2 of (Chowdhury et al., 2018), since  $A^{\top}S^{\top}\tilde{X}$  is the output for t=1 of their Algorithm 1. (Or rather, it is their output for each column of  $\tilde{X}$  and corresponding column of B.) To invoke their Theorem 2, we need to show that their equation (8) holds, which per their discussion following Theorem 3, holds for ridge leverage score sampling, with  $O(\varepsilon^{-2}\operatorname{sd}_{\lambda}\log\operatorname{sd}_{\lambda})$  samples, which our given  $m_R$  yields.

When we invoke their Theorem 2, we obtain

$$||A^{\top}S^{\top}\tilde{X} - X_1||_F \le \varepsilon(||X_1||_F + \frac{1}{\sqrt{\lambda}}||U_{k,\perp}B||_F).$$
(11)

Combining with (6) and using the triangle inequality, we have that, abbreviating  $\gamma_{\hat{A}_{(\lambda)},\hat{B}}$  as  $\gamma$ , up to the additive  $U_{k,\perp}$  term in (11), we have

$$\begin{split} \|A^{\top}S^{\top}\tilde{X} - X^*\|_F &\leq \|A^{\top}S^{\top}\tilde{X} - X_1\|_F + \|X_1 - X^*\|_F \\ &\leq \varepsilon \|X_1\|_F + \|X_1 - X^*\|_F + \frac{\varepsilon}{\sqrt{\lambda}} \|U_{k,\perp}B\|_F \\ &\leq \varepsilon \|X_*\|_F + (1+\varepsilon) \|X_1 - X^*\|_F + \frac{\varepsilon}{\sqrt{\lambda}} \|U_{k,\perp}B\|_F \\ &\leq \varepsilon \|X^*\|_F + 2\varepsilon\gamma \|X_*\|_F + \frac{\varepsilon}{\sqrt{\lambda}} \|U_{k,\perp}B\|_F \\ &\leq \varepsilon \left(1 + 2\gamma\right) \|X^*\|_F + \frac{\varepsilon}{\sqrt{\lambda}} \|U_{k,\perp}B\|_F \end{split}$$

for small enough  $\varepsilon$ , as claimed.

The time is dominated by that for computing  $\hat{A}^{-1}SB$ , where  $\hat{A}=SARR^{\top}A^{\top}S^{\top}$ , which we do via the conjugate gradient method. Via standard results (see, e.g., (Vishnoi, 2015), Thm 1.1), in  $O((T+m_S)\sqrt{\kappa(\hat{A})}\log(1/\alpha))d'$  time, where T is the time to compute the product of  $\hat{A}$  with a vector, we can obtain  $\tilde{X}$  with  $\|\tilde{X}-\hat{A}^{-1}SB\|_{\hat{A}}\leq \alpha\|\hat{A}^{-1}SB\|_{\hat{A}}$ , where the  $\hat{A}$ -norm is  $\|x\|_{\hat{A}}=x^{\top}\hat{A}x$ . Since S and R are (at least) constant-factor subspace embeddings, the singular values of SAR are within a constant factor of those of A, and so  $\kappa(\hat{A})$  is within a constant factor of

$$\kappa(AA^{\top} + \lambda I) = (\lambda + \sigma_1(A)^2) / (\lambda + \sigma_1(A)^2) = \kappa_{\lambda}^2.$$

We have

$$T = O(m_R m_S) = \tilde{O}(\varepsilon^{-2} \log m_S Z_{\lambda}^2 ||A||_F^2 \varepsilon^{-2} \hat{\kappa}^2 Z_{\lambda}^2 ||A||_F^2 \log(d))$$
  
=  $\tilde{O}(\varepsilon^{-4} \hat{\kappa}^2 Z_{\lambda}^4 ||A||_F^4 \log(d))$ 

Our running time is  $\tilde{O}(T\kappa_{\lambda}\log(1/\varepsilon))d'$ , with T as above. Translating to the notation using  $\psi$  terms, the result follows.  $\square$ 

#### C. Sampling from a Low-Rank Approximation

We need the following lemma, implied by the algorithm and analysis in Section 5.2 of (Boutsidis et al., 2016); for completeness we include a proof.

**Lemma C.1.** If  $S \in \mathbb{R}^{m_S \times n}$  and R are such that SA is a PCP of  $A \in \mathbb{R}^{n \times d}$ , and SAR is a PCP of SA, for error  $\varepsilon$  and rank t, and  $U \in \mathbb{R}^{m_S \times k}$  has orthonormal columns such that  $\|(I - UU^\top)SAR\|_F \le (1 + \varepsilon)\|SAR - [SAR]_t\|$ , then

$$Y^* = \operatorname{argmin}_Y \|YU^\top SA - A\|_F$$

has

$$||Y^*U^{\top}SA - A||_F \le (1 + O(\varepsilon))||A - A_t||_F.$$
 (12)

We also have

$$\|U^{\top}SA\|_F^2 \ge \|A_t\|_F^2 - O(\varepsilon)\|A\|_F^2.$$

*Proof.* Note that for matrix  $Y, Y(I - Y_t^+ Y_t) = (I - Y_t Y_t^+) Y$ , and that  $UU^\top SA$  is no closer to SA than is the projection of SA to the rowspace of  $U^\top SA$ , and that  $UU^\top = (SAR)_t(SAR)_t^+$  we have

$$\begin{split} \|A - Y^*U^\top SA\|_F &= \|A(I - (U^\top SA)^+ U^\top SA)\|_F \leq (1 + \varepsilon) \|SA(I - (U^\top SA)^+ U^\top SA)\|_F \\ &\leq (1 + \varepsilon) \|(I - UU^\top) SA\|_F \\ &\leq (1 + \varepsilon)^2 \|(I - UU^\top) SAR\|_F \\ &\leq (1 + \varepsilon)^3 \|(I - (SAR)_t (SAR)_t^+) SAR\|_F \\ &\leq (1 + \varepsilon)^3 \|(I - (SA)_t (SA)_t^+) SAR\|_F \\ &\leq (1 + \varepsilon)^4 \|(I - (SA)_t (SA)_t^+) SA\|_F \\ &\leq (1 + \varepsilon)^4 \|SA(I - (SA)_t^+ (SA)_t)\|_F \\ &\leq (1 + \varepsilon)^4 \|SA(I - A_t^+ A_t)\|_F \\ &\leq (1 + \varepsilon)^5 \|A(I - A_t^+ A_t)\|_F \\ &\leq (1 + \varepsilon)^5 \|A - A_t\|_F = (1 + O(\varepsilon)) \|A - A_t\|_F, \end{split}$$

as claimed.

For the last statement: we have  $\|SA\|_F^2 \geq (1-\varepsilon)\|A\|_F^2$ , since SA is a PCP, and by considering the projection of A onto the rowspans of blocks of t of its rows. We have also  $\|SA-[SA]_t\|_F^2 \leq (1+\varepsilon)\|A-[A]_t\|^2$ , using that SA is a PCP. Using these observations, we have

$$\begin{split} \|[SA]_t\|_F^2 &= \|SA\|_F^2 - \|SA - [SA]_t\|_F^2 \\ &\geq (1 - \varepsilon) \|A\|_F^2 - (1 + \varepsilon) \|A - [A]_t\|_F^2 \\ &= \|[A]_t\|_F^2 - \varepsilon (\|A\|_F^2 + \|A - [A]_t\|_F^2) \\ &\geq \|[A]_t\|_F^2 - 3\varepsilon \|A\|_F^2. \end{split}$$

Similarly,  $\|[SAR]_t\|_F^2 \ge \|[SA]_t\|_F^2 - 3\varepsilon \|SA\|_F^2$ , using that SAR is a PCP of SA. We then have, using these inequalities,

the PCP properties, and the hypothesis for U, that

$$\begin{split} \|U^{\top}SA\|_F^2 &= \|UU^{\top}SA\|_F^2 \\ &= \|SA\|_F^2 - \|(I - UU^{\top})SA\|_F^2 \\ &\geq (1 - \varepsilon)\|SAR\|_F^2 - (1 + \varepsilon)^2\|SAR - [SAR]_t\|_F^2 \\ &\geq \|[SAR]_t\|_F^2 - 4\varepsilon\|SAR\|_F^2 \\ &\geq (\|[SA]_t\|_F^2 - 3\varepsilon\|SA\|_F^2) - 4(1 + \varepsilon)\varepsilon\|SA\|_F^2 \\ &\geq (\|[A]_t\|_F^2 - 3\varepsilon\|A\|_F^2) - 3\varepsilon(1 + \varepsilon)\|A\|_F^2 - 4(1 + \varepsilon)^2\varepsilon\|A\|_F^2 \\ &\geq \|[A]_t\|_F^2 - 13\varepsilon\|A\|_F^2, \end{split}$$

for small enough  $\varepsilon$ , and the last statement of the lemma follows.

Before a proof, we give a re-statement of Theorem 1.5.

**Theorem C.2** (Dynamic Data Structure for LRA). Given a matrix  $A \in \mathbb{R}^{n \times d}$ , target rank k, and estimate  $\hat{\sigma}_k \leq \sigma_k(A)$ , error parameter  $\varepsilon > 0$ , and estimate  $\tau$  of  $\|A - A_k\|_F^2$ , there exists a data structure representing a matrix  $Z \in \mathbb{R}^{n \times d}$  with rank k such that if  $\|A_k\|_F^2 \geq \varepsilon \|A\|_F^2$ , with probability at least 99/100,  $\|A - Z\|_F^2 \leq (1 + O(\varepsilon)) \|A - A_k\|_F^2$ , where  $A_k$  is the best rank-k approximation to A. Further, the time taken to construct the representation of Z is

$$\tilde{O}(\varepsilon^{-6}k^3 + \varepsilon^{-4}\psi_{\lambda}(\psi_{\lambda} + k^2 + k\psi_k)),$$

where  $\psi_{\lambda} = \|A\|_F^2/(\tau/k + \hat{\sigma}_k^2)$  and  $\psi_k = \|A\|_F^2/\sigma_k(A)$ . Given  $j \in [d]$ ,  $i \in [n]$  can be generated with probability  $(Z)_{ij}^2/\|Z)_{*,j}\|^2$  in expected time  $O(\|A\|_F^2/\hat{\sigma}_k^2 + m_R^2\kappa^2)$ , where  $\kappa$  is the condition number of A, and  $m_R = O(k \log k + \varepsilon^{-1}k)$ .

*Proof.* The matrix Z is the implicit output of Algorithm 3. In that algorithm, the choice of  $m_S = O(\hat{m}_S Z_{\lambda}^2 \|A\|_F^2)$  rows constitutes an effective  $k\hat{m}_S = O(\varepsilon^{-2}k\log k)$  ridge-leverage score samples of the rows of A. We assume that the input  $\tau$  is within a constant factor of  $\|A - A_k\|_F^2$ , so that  $\lambda = \tau/k$  is within a constant factor of  $\|A - A_k\|_F^2/k$ . Theorem 6 of (Cohen et al., 2017) implies that under these conditions, SA will be a rank-k Projection-Cost Preserving (PCP) sketch of A with error parameter  $\varepsilon$ , a  $(k, \varepsilon)$ -PCP.

Similarly to S,  $R_1$  will be a (column) rank-k PCP of SA, here using that the PCP properties of SA imply that  $\|(S(A-A_k)\|_F^2=(1\pm\varepsilon)\|A-A_k\|_F^2$ , and so the appropriate  $\lambda$ , and  $Z_\lambda$ , for SA are within constant factors of those for A. Let  $\hat{A}=SAR_1$ . Lemma 16 and Theorem 1 of (Cohen et al., 2017) imply that applying their Algorithm 1 to  $\hat{A}$  yields  $S_2\in\mathbb{R}^{m_{S_2}\times m_S}$  so that  $S_2\hat{A}$  is a  $(k,\varepsilon)$ -PCP for  $\hat{A}$ , and similarly  $S_2\hat{A}R_2$  is a  $(k,\varepsilon)$ -PCP for  $S_2\hat{A}$ .

We apply Lemma C.1 with  $\hat{A}^{\top}$ ,  $R_2^{\top}$ ,  $S_2^{\top}$ , and  $V^{\top}$  in the roles of A, S, R, and U in the lemma. We obtain that  $\tilde{Y}=(\hat{A}R_2V)^+\hat{A}=\operatorname{argmin}_Y\|\hat{A}R_2VY-\hat{A}\|_F$  has  $\|\hat{A}R_2V\tilde{Y}-\hat{A}\|_F\leq (1+O(\varepsilon))\|\hat{A}-\hat{A}_k\|_F$ , that is, U as constructed in Algorithm 3 has  $UU^{\top}\hat{A}=\hat{A}R_2V\tilde{Y}$ , and therefore satisfies the conditions of Lemma C.1 for A, S,  $R_1$ . This implies that  $Y^*=A(U^{\top}SA)^+=\operatorname{argmin}_Y\|YU^{\top}SA-A\|_F$  has

$$||Y^*U^{\top}SA - A||_F \le (1 + O(\varepsilon))||A - A_k||_F.$$
 (13)

It remains to solve the multiple-response regression problem  $min_Y ||YU^\top SA - A||_F$ , which we do more quickly using the samplers  $R_3$  and  $R_4$ .

We next show that  $R_3^{\top}$  is a subspace  $\varepsilon_0$ -embedding of  $(U^{\top}SA)^{\top}$ , and supports low-error matrix product estimation, so that Thm. 36 of (Clarkson & Woodruff, 2013) can be applied. Per Lemma 3.1 and per Lemma 32 of (Clarkson & Woodruff, 2013),  $k\hat{m}_{R_3} = O(\varepsilon_0^{-2}k\log k + \varepsilon^{-1}k)$  leverage-score samples of the columns of  $U^{\top}SA$  suffice for these conditions to hold

To obtain  $k\hat{m}_{R_3}$  leverage score samples, we show that  $1/\varepsilon$  length-squared samples of the columns of SA suffice to contain one length-squared sample of  $U^{\top}SA$ , and also that  $\|(U^{\top}SA)^{+}\| \leq 1/\hat{\sigma}_{k}$ , using the input condition on  $\hat{\sigma}_{k}$  that  $\sigma_{k}(A) \geq \hat{\sigma}_{k}$ , so that the given value of  $m_{R_3}$  in the call to Lensquared for  $R_3$  is valid.

For the first claim, by hypothesis  $\|A_k\|_F^2 \ge \varepsilon \|A\|_F^2$ , and by adjusting constants, U as computed satisfies the conditions of Lemma C.1 for some  $\varepsilon' = \alpha \varepsilon$  for constant  $\alpha > 0$ , so by that lemma and by hypothesis

$$\|U^{\top}SA\|_F^2 \ge \|A_k\| - O(\alpha\varepsilon)\|A\|_F^2$$

$$\ge \varepsilon(1 - O(\alpha))\|A\|_F^2$$

$$\ge \varepsilon(1 - O(\alpha))(1 - \varepsilon)\|SA\|_F^2,$$

so adjusting constants, we have  $\|U^{\top}SA\|_F^2 \geq \varepsilon \|SA\|_F^2$ . Using that U has orthonormal columns, we have for  $j \in [d]$  that  $\|U^{\top}SA_{*,j}\|/\|U^{\top}SA\|_F^2 \leq \|SA_{*,j}\|/\varepsilon \|SA\|_F^2$ , so the probability of sampling j using length-squared probabilities for SA is least  $\varepsilon$  times that for  $U^{\top}SA$ .

For the claim for the value of  $m_{R_3}$  used for  $R_3$ , using the PCP properties of SA and  $SAR_1$ , we have

$$\sigma_k(U^{\top}SA) = \sigma_k(UU^{\top}SAR_1) = \sigma_k(SAR_1) \ge (1 - \varepsilon)\sigma_k(SA) \ge (1 - O(\varepsilon))\sigma_k(A).$$

so the number of length-squared samples returned by LENSQSAMPLE suffice.

So using Thm. 36 of (Clarkson & Woodruff, 2013),  $\tilde{Y}_3 = \operatorname{argmin}_Y \|(YU^{\top}SA - A)R_3\|_F$  satisfies

$$\|\tilde{Y}_3 U^\top SA - A\|_F \leq (1+\varepsilon) \min_Y \|Y U^\top SA - A\|_F \leq (1+O(\varepsilon)) \|A - A_k\|_F,$$

where the last inequality follows from (12).

Similar conditions and results can be applied to direct leverage-score sampling of the columns of  $U^{\top}SAR_3$ , resulting in  $\tilde{Y}_4 = \min_Y \|(YU^{\top}SAR_3 - AR_3)R_4\|_F$ , where there is  $m_{R_4} = O(\varepsilon_0^{-2}k\log k + \varepsilon^{-1}k)$  such that these conditions hold for  $R_4$ . This implies  $\tilde{Y}_4$  is an approximate solution to  $\min_Y \|(YU^{\top}SA - A)R_3\|_F$ , and therefore  $AR\tilde{Y}_4 = AR(U^{\top}SAR)^+$  has  $\|AR\tilde{Y}_4U^{\top}SA - A\|_F \le (1+O(\varepsilon))\|A - A_k\|_F$ , as claimed. We have  $W \leftarrow (U^{\top}SAR)^+U^{\top} = \tilde{Y}_4U^{\top}$ , so the claimed output condition on ARWSA holds.

Turning to the time needed, Lemma 16 and Theorem 1 of (Cohen et al., 2017) imply that the time needed to construct  $S_2$  and  $R_2$  is

$$O(m_{R_1}m_S + k^2m_S) = O(m_S(m_S + k^2))$$
(14)

The time needed to construct V from  $S_2SAR_1R_2 \in \mathbb{R}^{m_{S_2} \times m_{S_2}}$  is

$$O(m_{S_2}^3) = \tilde{O}(\varepsilon^{-6}k^3) \tag{15}$$

The time needed to construct U from  $V \in \mathbb{R}^{m_{R_2} \times k}$  and  $SAR_1R_2 \in \mathbb{R}^{m_S \times m_{R_2}}$  by multiplication and QR factorization is

$$O(km_{R_2}m_S + k^2m_S) = \tilde{O}(m_Sk^2\varepsilon^{-2})$$
(16)

Computation of  $U^{\top}SAR_3$  requires  $O(km_Sm_{R_3})$  time, where  $m_{R_3} = O(\hat{m}_{R_3}\varepsilon^{-1}Z_k^2\|A\|_F^2$ , and  $\hat{m}_{R_3} = O(\log k + \varepsilon^{-1})$ , that is,

$$\tilde{O}(m_S k \varepsilon^{-2} Z_k^2 ||A||_F^2) \tag{17}$$

time. Leverage-score sampling of the rows of  $(U^{\top}SAR_3)^{\top} \in \mathbb{R}^{m_{R_3} \times k}$  takes, applying Theorem A.4 and using  $m_{R_4} = O(k(\log k + \varepsilon^{-1}))$ , time at most

$$O(\log(m_{R_3})km_{R_3} + k^{\omega}\log\log(km_{R_3}) + k^2\log m_{R_3} + m_{R_4}km_{R_3}^{1/\log(m_{R_3})})$$
(18)

$$= \tilde{O}(k\varepsilon^{-2}Z_k^2||A||_F^2 + k^\omega + k^2\varepsilon^{-1})$$
(19)

Computation of  $(U^{\top}SAR)^+$  from  $U^{\top}SAR$  requires  $O(k^2m_{R_4}) = \tilde{O}(\varepsilon^{-1}k^3)$  time. (With notation that R has  $m_R = m_{R_4}$  columns.) Given  $(U^{\top}SAR)^+$ , computation of  $W = (U^{\top}SAR)^+U^{\top}$  requires

$$O(km_R m_S) = \tilde{O}(m_S k^2 \varepsilon^{-1}) \tag{20}$$

time. Putting together (14),(16), (17), (20), (15), (18), we have

$$O(m_{S}(m_{S} + k^{2})) + \tilde{O}(m_{S}k^{2}\varepsilon^{-2}) + \tilde{O}(m_{S}k\varepsilon^{-2}Z_{k}^{2}||A||_{F}^{2}) + \tilde{O}(m_{S}k^{2}\varepsilon^{-1})$$

$$+ \tilde{O}(\varepsilon^{-6}k^{3}) + tO(k\varepsilon^{-2}Z_{k}^{2}||A||_{F}^{2} + k^{\omega} + k^{2}\varepsilon^{-1})$$

$$= \tilde{O}(m_{S}(m_{S} + k^{2}\varepsilon^{-2} + k\varepsilon^{-2}Z_{k}^{2}||A||_{F}^{2}) + \varepsilon^{-6}k^{3})$$

Here  $m_S = O(\hat{m}_S Z_{\lambda}^2 ||A||_F^2) = \tilde{O}(\varepsilon^{-2} Z_{\lambda}^2 ||A||_F^2)$ , so the time is

$$\tilde{O}(\varepsilon^{-6}k^3 + \varepsilon^{-4}Z_{\lambda}^2 \|A\|_F^2 (Z_{\lambda}^2 \|A\|_F^2 + k^2 + kZ_k^2 \|A\|_F^2))$$

Queries as in the theorem statement for given  $j \in [d]$  can be answered as in (Tang, 2019), in the time given. Briefly: given j, let  $v \leftarrow (WSA)_{*,j} \in \mathbb{R}^{m_R}$ . Let  $\hat{A}$  denote AR. Using LENSQSAMPLE (and large  $m_R$ ), generate a sampling matrix  $S_3$  with  $O(Z^2\|\hat{A}\|^2)$  rows, and estimate  $\|\hat{A}v\|^2 \approx \beta_v \leftarrow \|S\hat{A}v\|^2$ . Generate  $i \in [n]$  via rejection sampling as follows. In a given trial, pick  $j^* \in [d]$  with probability proportional to  $\|\hat{A}_{*,j}\|^2 v_j^2$  using DYNSAMP(A), then pick  $i \in [m_R]$  with probability  $\hat{A}_{i,j^*}^2/\|\hat{A}_{*,j^*}\|^2$ . This implies that  $i \in [n]$  has been picked with probability  $p_i = \sum_j \hat{A}_{ij}^2 v_j^2/\sum_j \|\hat{A}_{*,j}\|^2 v_j^2$ . Now for a value  $\alpha > 0$ , accept with probability  $q_i/\alpha p_i$ , where  $q_i = (\hat{A}_{i,*}v)^2/\beta_v$ , otherwise reject. This requires  $\alpha \geq q_i/p_i$ . and takes expected trials  $\alpha$ . So an upper bound is needed for

$$\frac{q_i}{p_i} = \frac{(\hat{A}_{i,*}v)^2}{\beta_v} \frac{\sum_j \|\hat{A}_{*,j}\|^2 v_j^2}{\sum_j \hat{A}_{ij}^2 v_j^2}.$$

We have  $(\hat{A}_{i,*}v)^2 \leq m_R \sum_j \hat{A}_{ij}^2 v_j^2$  using Cauchy-Schwarz, and  $\beta_v \approx \|\hat{A}v\|^2 \geq \|v\|^2/\|\hat{A}^+\|^2$ , and also  $\sum_j \|\hat{A}_{*,j}\|^2 v_j^2 \leq \|v\|^2 \max_j \|\hat{A}_{*,j}\|^2 \leq \|v\|^2 \|\hat{A}\|^2$ . Putting this together  $\alpha \geq m_R \|\hat{A}\|^2 \|\hat{A}^+\|^2 = O(m_R \kappa(A))$  will do. The work per trial is  $O(m_R \log(nd))$ , and putting that together with the time to compute  $\beta_v$ , the theorem follows.  $\square$ 

#### **D.** Additional Experiments

#### D.1. Low-Rank Approximation

As we stated in Section 6.1, although the accuracy of our algorithm is slightly worse, by increasing the sample size slightly, our algorithm achieves a similar accuracy as (Arrazola et al., 2020), but still has a faster runtime. The results are shown in Table 5. Here we set (r,c)=(500,800) for MovieLens 100K and (r,c)=(700,1100) for KOS data for our algorithms, but do not change (r,c) for the algorithm in (Arrazola et al., 2020) as in Section 6.1.

Table 5. Performance of our algorithm and ADBL on MovieLens 100K and KOS data, respectively.							
	k = 10	k = 15	k = 20		k = 10	k = 15	k = 20
$\varepsilon$ (Ours)	0.0323	0.0439	0.0521	$\varepsilon(\text{Ours})$	0.0291	0.0390	0.0476
$\varepsilon(ADBL)$	0.0262	0.0424	0.0538	$\varepsilon(ADBL)$	0.0186	0.0295	0.0350
Runtime (Ours, Query)	0.341s	0.365s	0.370s	Runtime (Ours, Query)	0.826s	0.832s	0.831s
Runtime				Runtime			
(Ours, Total)	0.412s	0.417s	0.415s	(Ours, Total)	0.979s	0.994s	0.986s
Runtime (ADBL, Query)	0.863s	0.917s	1.024s	Runtime (ADBL, Query)	1.501s	1.643s	1.580s
Runtime (ADBL, Total)	0.968s	1.003s	1.099s	Runtime (ADBL, Total)	1.814s	1.958s	1.897s
Runtime of SVD 2.500s		Runtime of SVD		36.738s			

Table 5. Performance of our algorithm and ADBL on MovieLens 100K and KOS data, respectively.