

Data Driven Adaptive Sampling for Low-Rank Matrix Approximation

Arvind Rathnashyam*

RPI Math and CS, rathna@rpi.edu

Nicolas Boullé

Cambridge Math, nb690@cam.ac.uk

Alex Townsend

Cornell Math, ajt453@cornell.edu

Abstract

We consider the problem of low-rank matrix approximation in case the when the matrix \mathbf{A} is accessible only via matrix-vector products and we are given a budget of $k + p$ matrix-vector products. This situation arises in practice when the cost of data acquisition is high, despite the Numerical Linear Algebra (NLA) costs being low. We create an adaptive sampling algorithm to optimally choose vectors to sample. The Randomized Singular Value Decomposition (rSVD) is an effective algorithm for obtaining the low rank representation of a matrix developed by [17]. Recently, [2] generalized the rSVD to Hilbert-Schmidt Operators where functions are sampled from non-standard Covariance Matrices when there is already prior information on the right singular vectors within the column space of the target matrix, \mathbf{A} . In this work, we develop an adaptive sampling framework for the Matrix-Vector Product Model which does not need prior information on the matrix \mathbf{A} . We provide a novel theoretic analysis of our algorithm with subspace perturbation theory. We extend the analysis of [29] for right singular vector approximations from the randomized SVD in the context of non-symmetric rectangular matrices. We also test our algorithm on various synthetic, real-world application, and image matrices. Furthermore, we show our theory bounds on matrices are stronger than state-of-the-art methods with the same number of matrix-vector product queries.

*Work was partially completed at Cornell Summer 2023

1 Introduction

In many real-world applications, it is often not possible to run experiments in parallel. Consider the following setting, there are a set of n inputs and m outputs, and there exists a PDE such it maps any set of inputs in $\mathbb{C}^m \rightarrow \mathbb{C}^n$. However, to run experiments, it takes hours for set up, execution, or it is expensive, e.g. aerodynamics [14], fluid dynamics [21]. Thus, after each experimental run, we want to sample a function such that in expectation, we will be exploring an area of the PDE which we have the least knowledge of. For Low-Rank Approximation the Randomized SVD, [17], has been theoretically analyzed and used in various applications. Even more recently, [4] discovered if we have prior information on the right singular vectors of \mathbf{A} , we can modify the Covariance Matrix such that the sampled vectors are within the column space of \mathbf{A} . They extended the theory for Randomized SVD where the covariance matrix is now a general PSD matrix. The basis of our analysis is the idea of sampling vectors in the Null-Space of the Low-Rank Approximation. This idea has been introduced recently in Machine Learning in [31] for training neural networks for sequential tasks. In a Bayesian sense, we want to maximize the expected information gain of the PDE in each iteration by sampling in the space where we have no information. This leads to the formulation of our iterative algorithm for sampling vectors for the Low-Rank Approximation. The current state of the art algorithms for low-rank matrix approximation in the matrix-vector product model used a fixed covariance matrix structure. In this paper, we consider the adaptive setting where the algorithm \mathcal{A} chooses a vector $\mathbf{v}^{(k)}$ with access to the previous query vectors $\mathbf{v}_1, \dots, \mathbf{v}^{(k-1)}$, the matrix-vector products $\mathbf{A}\mathbf{v}^{(1)}, \dots, \mathbf{A}\mathbf{v}^{(k-1)}$, and the intermediate low-rank matrix approximations, $\mathbf{Q}^{(k)}(\mathbf{Q}^{(k)})^H \mathbf{A}$, where $\mathbf{Q}^{(k)} \triangleq \text{orth}(\mathbf{A}\mathbf{V}^{(k)})$ where $\mathbf{V}^{(k)}$ is the concatenation of vectors $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k)}$ and $\mathbf{Q}^{(k)} \triangleq \text{orth}(\mathbf{V}^{(k)})$.

Adaptive Sampling techniques for Low-Rank Matrix Approximation first appeared in CUR Matrix Decomposition in [15]. Optimal column-sampling for the CUR Matrix Decomposition received much attention as can be seen in the works [18, 10, 9]. More recently, [27] gave an algorithm for sampling the rows for CUR-Matrix Factorization and proved error bounds by induction. Similar to adaptively choosing a function, in recommender systems, the company can ask users for surveys and obtain data with high probability is a better representation of the column space of \mathbf{A} than a random sample. Choosing the right people to give an incentivized survey (e.g. gift card upon completion) can save a company significant expenses.

The theoretical properties of adaptively sampled matrix vector queries have been studied in [6]. Their bounds are used in [1] to develop adaptive bounds for their low-rank matrix approximation method using Krylov Subspaces. To our knowledge, we are the first paper to give an algorithm for low-rank approximation in the non-symmetric matrix low-rank approximation in the matrix-vector product model. Our algorithm utilizes the SVD computation of the low-rank approximation at each step to sample the next vector. Although there are runtime limitations, both in theory under certain conditions and most real-world matrices, our algorithm gets the most value out of each sampled vector.

We will now clearly state our contributions.

Main Contributions.

1. We develop a novel adaptive sampling algorithm for Low-Rank Matrix Approximation problem in the matrix-vector product model which does not utilize prior information of \mathbf{A} .
2. We provide a novel theoretical analysis which utilizes subspace perturbation theory.
3. We perform extensive experiments on matrices with various spectrums and compare with the state of the art methods.

2 Notation, Background Materials, and Relevant Work

In this section we will introduce the notation we use throughout the paper, perturbations of singular spaces, as well as relevant work in the Low-Rank Matrix Approximation Literature.

2.1 Notation

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ represent the target matrix. $\|\cdot\|$ represents the spectral norm, which is equivalent to the max eigenvalue of the argument, $\sigma_{\max}(\cdot)$. Quasimatrices (matrices with infinite rows and finite columns) will be denoted as a variation of the symbol, $\mathbf{\Omega}$. The pseudoinverse is represented by $(\cdot)^\dagger$ s.t. $\mathbf{X}^\dagger = (\mathbf{X}^H \mathbf{X})^{-1} \mathbf{X}^H$. The Projection Matrix is defined as $\Pi_{\mathbf{Y}} = \mathbf{Y} \mathbf{Y}^\dagger = \mathbf{Y} (\mathbf{Y}^H \mathbf{Y})^{-1} \mathbf{Y}^H$ as the projection on to the column space of \mathbf{Y} . If \mathbf{Y} has orthogonal columns, then $\Pi_{\mathbf{Y}}$ is the Orthogonal Projection defined as $\Pi_{\mathbf{Y}} = \mathbf{Y} \mathbf{Y}^H$. Let $a \wedge b = \min(a, b)$ and $a \vee b = \max(a, b)$. Let $\mathbb{O}_{n,k}$ be the set of all $n \times k$ matrices with orthogonal columns, i.e. $\{\mathbf{V} : \mathbf{V}^H \mathbf{V} = \mathbf{I}_{k \times k}\}$. We also denote $\mathcal{MN}(\mathbf{0}, \mathbf{I}_{n \times n}, \mathbf{I}_{m \times m})$, denote the distribution of $m \times n$ standard gaussian matrices. The Frobenius norm for a matrix is defined as,

$$\|\mathbf{A}\|_F = \left(\sum_{i \in [m]} \sum_{j \in [n]} \mathbf{A}_{i,j}^2 \right)^{1/2} = \sqrt{\text{Tr}(\mathbf{A}^H \mathbf{A})} = \sqrt{\text{Tr}(\mathbf{A} \mathbf{A}^H)} \quad (1)$$

We define $[\mathbf{A}]_r$ as the best rank- r approximation to \mathbf{A} w.r.t the Frobenius norm. We use Big-O notation, $y \leq O(x)$, to denote $y \leq Cx$ for some positive constant, C . We define \mathbb{E} as expectation, \mathbb{P} as probability, and \mathbb{V} as variance. We will denote normal text characters A, X, Y as matrices and lower roman boldface characters $\mathbf{x}, \mathbf{y}, \mathbf{z}$ as vectors.

2.2 Singular Subspace Perturbations

To represent the distance between subspaces we utilize the $\sin \Theta$ norm. Let \mathcal{X}, \mathcal{Y} be subspaces, then we denote the principal angles between subspaces (PABS) \mathcal{X} and \mathcal{Y} as $\frac{\pi}{2} \geq \Theta_1(\mathcal{X}, \mathcal{Y}) \geq \dots \geq \Theta_{m \wedge n}(\mathcal{X}, \mathcal{Y})$.

2.3 Relevant Works

The Randomized Singular Value Decomposition was developed and analyzed thoroughly in [17]; throughout this paper we will refer to this algorithm as HMT. The review work by [23] gives significant theory on the Randomized SVD. [2] proposed learning the Hilbert-Schmidt Operators associated with the Green's Functions with the randomized SVD algorithm. One of their key findings is they can better approximate the HS Operator when they use functions drawn from $\mathcal{GP}(\mathbf{0}, \mathbf{K})$ where \mathbf{K} is not the identity. [3] extended upon previous work on generalizing the Randomized SVD to learning HS Operators.

The most relevant work to ours is likely [12]. The measure of accuracy in the Krylov Subspace is measured by the $\sin \Theta$ norm. We would like to note the Krylov Subspace method takes q times more matrix-vector products and thus is not a suitable method for our problem. Work similar to ours with regards to upper bounding the sine of the principal angles between subspaces in the context of the Randomized SVD is explored in [11] and [28].

A similar analysis of a power method is explored in [19] utilizing subspace perturbation theory. In this work, they consider the Matrix-Vector products have noise. In this work, similarly to [12], it takes d times more matrix-vector products to recover the right singular space. Furthermore, a

similar projection-based analysis based on the sines of the singular vector perturbations is done in [22].

3 Data Driven Sampling

In this section, we will go over the covariance matrices proposed papers and we consider choosing the optimal covariance matrix adaptively for sampling vectors. In the seminal paper by [17], the covariance matrix is given as identity matrix, $\mathbf{C} \triangleq \mathbf{I}$. In the generalization of the Randomized SVD, when given some prior information of the matrix, the covariance matrix is given as $\mathbf{C} \triangleq \mathbf{K}$ where \mathbf{K} has some information on the right singular vectors of \mathbf{A} (e.g. discretization of Green's Function of a PDE). Let $\tilde{\mathbf{V}}$ be the right singular vectors of the SVD of the low-rank approximation at iteration $k - 1$, then the update for the covariance matrix is given as $\mathbf{C}^{(k+1)} \triangleq \tilde{\mathbf{V}}_{(:,k)} \tilde{\mathbf{V}}_{(:,k)}^H$. Throughout this paper we will only consider $\mathbf{C}^{(0)} = \mathbf{I}$, however using theory from [2], this can be extended to $\mathbf{C}^{(0)} = \mathbf{K}$ if one has some knowledge of the right singular vectors. A similar algorithm can be found in [31]. Naturally, want to continuously sample in the null space of the the matrix approximation we have already obtained. This ensures we are learning new information in each iteration as we don't want to 'waste' samples which do not learn any new information about the matrix. To further motivate our covariance update, we will introduce the following remark.

Remark 1. Let $\mathbf{U}\Sigma\mathbf{V}^H$ be the SVD of \mathbf{A} , then the Covariance update described in Section 3 is the optimal covariance update is the optimal covariance matrix for sampling vectors at iteration k .

Remark 1 is an intuitive result, in that when we are learning a matrix \mathbf{A} , we would optimally want to sample the right singular vectors, so the resultant matrix product is the left singular vectors.

3.1 Algorithm

The Pseudo Code for the optimal function sampling is given in Algorithm 1. For efficient updates, we frame all operations as rank-1 updates.

In Algorithm 1, we first sample a standard normal gaussian matrix which can be considered as the oversampling vectors. These oversampling vectors are used to approximate the first singular vector. This is the first vector which is *adaptively* sampled. Next, we form the low-rank approximation $\mathbf{Q}\mathbf{Q}^H\mathbf{A}$ with the adaptive matrix vector query. From here, we adaptively query with the k th right singular value of the SVD of the the low rank approximation at iteration $k - 1$. We believe this algorithm to be the closest to replicate the idea given in Remark 1.

4 Theory

In this section we will give the mathematical setup for the theoretical analysis. We will then represent theorems from relevant works on the error bounds for their low-rank approximation methods. We will then give our error bounds and general theory of Algorithm 1 with the proofs in the appendix.

4.1 Setup.

We follow a similar setup as previous literature. Let $\rho \triangleq \text{rank}(\mathbf{A}) \leq m \wedge n$, we will factorize \mathbf{A} as

$$\mathbf{A} = \begin{bmatrix} \mathbf{U}_k & \mathbf{U}_{\rho-k} \end{bmatrix} \begin{bmatrix} \Sigma_k & \\ & \Sigma_{\rho-k} \end{bmatrix} \begin{bmatrix} \mathbf{V}_k^H \\ \mathbf{V}_{\rho-k}^H \end{bmatrix} = \sum_{i=1}^{\rho} \sigma_i \mathbf{u}_i \mathbf{v}_i^H = \sum_{i=1}^{\rho} \mathbf{U}_{(i)} \Sigma_{(i)} \mathbf{V}_{(i)}^H \quad (2)$$

Algorithm 1 Bayesian Function Sampling

1: **Input:** HS Operator: \mathcal{F} , Rank: r , Initial Covariance: \mathbf{C} , Oversampling Parameter: p
2: **Output:** Rank- r Approximation, $\hat{\mathbf{A}}_r$
3: $\mathbf{\Omega} \leftarrow \underbrace{[\mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \overset{\text{i.i.d.}}{\vdots} \quad \mathcal{N}(\mathbf{0}, \mathbf{C})]}_p \triangleright$ Sample Oversampling Vectors from Standard Normal Matrix
4: $\mathbf{Y} \leftarrow \mathbf{A}\mathbf{\Omega}$ \triangleright Matrix Vector Products
5: $[\mathbf{Q}_0, \sim] \leftarrow \text{QR}(\mathbf{Y})$ \triangleright Find Orthonormal Basis
6: $\tilde{\mathbf{A}} \leftarrow \mathbf{0}_{m \times n}$ \triangleright Initial Low-Rank Approximation
7: **for** $k \in 1, 2, \dots, r$ **do**
8: $\tilde{\mathbf{A}}_k \leftarrow \tilde{\mathbf{A}}_{k-1} + \mathbf{Q}_{(k-1)} \mathbf{Q}_{(k-1)}^H \mathbf{A}$ \triangleright Rank-1 update to the low-rank approximation
9: $[\tilde{\mathbf{U}}, \tilde{\mathbf{\Sigma}}, \tilde{\mathbf{V}}] \leftarrow \text{SVD}(\tilde{\mathbf{A}}_k)$ \triangleright SVD of current low-rank approximation
10: $\mathbf{C}^{(k+1)} \leftarrow \tilde{\mathbf{V}}_{(:,k)} \tilde{\mathbf{V}}_{(:,k)}^H$ \triangleright Form new Covariance Matrix
11: $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}^{(k+1)})$ \triangleright Adaptive Sampling of Vector
12: $\mathbf{Y} \leftarrow [\mathbf{Y} \quad \mathbf{A}\mathbf{x}]$ \triangleright Matrix-Vector Product
13: $[\mathbf{Q}_k, \sim] \leftarrow \text{QR}(\mathbf{Y})$ \triangleright Find Orthonormal Basis
14: **end for**
15: $\tilde{\mathbf{A}}_r \leftarrow \mathbf{Q}_r \mathbf{Q}_r^H \mathbf{A}$ \triangleright Final Low-Rank Approximation
16: **return:** $\tilde{\mathbf{A}}_r$

Furthermore, we let $\mathbf{A}_{(k)} \triangleq \sigma_k \mathbf{u}_k \mathbf{v}_k^H$. Let $\mathbf{\Omega} \in \mathbb{R}^{n \times \ell}$ be a test matrix where $\ell = k + p$ denotes the number of samples and p is the oversampling parameter.

4.2 Query Lower Bound for Frobenius Norm

In this section, we give information theoretic lower bounds on query complexity. We assume the algorithm, \mathcal{A} not only has access to the matrix-vector products, but also has available the SVD of the intermediate low-rank approximations.

Theorem 2. *There exists a distribution \mathcal{D} over $\mathbb{C}^{m \times n}$ such that for $\mathbf{A} \sim \mathcal{D}$ there exists an adaptive algorithm (possibly randomized) with access to vector queries $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(k-1)}$ where w.l.o.g $\|\mathbf{v}^{(i)}\| = 1$ for all $i \in [k-1]$ and $(\mathbf{v}^{(i)})^H \mathbf{v}^{(j)} = \delta_{ij}$ for all $i, j \in [k-1]$, matrix-vector queries, $\mathbf{A}\mathbf{v}^{(1)}, \dots, \mathbf{A}\mathbf{v}^{(k-1)}$, and intermediate low-rank approximations, $\mathbf{Q}^{(1)}(\mathbf{Q}^{(1)})^H \mathbf{A}, \dots, \mathbf{Q}^{(k-1)}(\mathbf{Q}^{(k-1)})^H \mathbf{A}$, which requires $k = O(\Xi)$ vector queries to obtain a rank- k matrix with orthogonal columns, \mathbf{Q} , such that with probability at least $1 - \delta$ for a $\delta \in (0, 1)$ such that*

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H \mathbf{A}\|_F \leq (1 + \varepsilon) \min_{\substack{\mathbf{U} \in \mathbb{C}^{m \times k} \\ \mathbf{U}^H \mathbf{U} = \mathbf{I}_k}} \|\mathbf{A} - \mathbf{U}\mathbf{U}^H \mathbf{A}\|_F \quad (3)$$

Proof.

4.3 Analysis of Algorithm 1

First we will introduce a lemma for the resultant vector of sampling from $\mathbf{C}^{(k)}$. Since our general proof technique will be an induction. We first want to understand how well we are able to approximate the first right singular vector. To do this, we must know the singular vector perturbation from the error of the low-rank matrix approximation.

Lemma 3. Let $\mathbf{A} \in \mathbb{C}^{m \times n}$ and \mathbf{Q} be an orthogonal matrix representing the basis of the subspace of $\mathbf{Y} \in \mathbb{C}^{m \times k}$. Let $\mathbf{v} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$ and $\tilde{\mathbf{v}}$ represent the k -th right singular vector of $\mathbf{Q}\mathbf{Q}^H\mathbf{A}$. Denote $\hat{\mathbf{Q}} \triangleq \text{orth}([\mathbf{Y} \ \mathbf{v}])$ and $\tilde{\mathbf{Q}} \triangleq \text{orth}([\mathbf{Y} \ \tilde{\mathbf{v}}])$, if $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H\mathbf{A}\|_F \leq C\sigma_{k+1}$ and $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H\mathbf{A}\|_2 \leq c\sigma_{k+1}$ for positive constants $c, C > 0$ where $C \geq c$, then we have with probability $1 - \delta$

Claim (i):

$$\|\mathbf{A} - \tilde{\mathbf{Q}}\tilde{\mathbf{Q}}^H\mathbf{A}\|_F \leq C\sigma_{k+1} \sqrt{1 - \frac{(\sigma_k - 1)^2}{C\sigma_{k+1}^2} + \frac{2c}{C} (1 + \sigma_k^{-2})} \quad (4)$$

Claim (ii):

$$\|\mathbf{A} - \hat{\mathbf{Q}}\hat{\mathbf{Q}}^H\mathbf{A}\|_F \leq C\sigma_{k+1} \sqrt{1 - O\left(\frac{1}{r}\right) \left(1 - \sqrt{\frac{1}{c_1} \log \frac{2}{\delta}}\right)} \quad (5)$$

Proof. The proof is deferred to Appendix A.1. ■

It is clear to see in Equation (4) the strength of Bayesian sampling described in Algorithm 1 when there is sufficient singular value gap. When σ_k/σ_{k+1} is sufficiently large with respect to σ_k , then it follows that the $(\sigma_k - 1)^2/(C\sigma_{k+1}^2)$ term will dominate the $\Omega(1 + \sigma_k^{-2})$ term.

Theorem 4 (Sufficient Singular Value Gap). If $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H\mathbf{A}\|_F \leq C\sigma_{k+1}$ and $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H\mathbf{A}\|_2 \leq c\sigma_{k+1}$ for positive constants $c, C > 0$ where $C \geq c$, then we have with probability $1 - \delta$ Bayesian sampling described in Section 3 decreases the low-rank approximation error faster than a normal vector sample when

$$\frac{\sigma_k}{\sigma_{k+1}} \geq \sqrt{\frac{C^2}{r} \left(1 - \sqrt{\frac{1}{c_1} \log \frac{2}{\delta}}\right) + 2Cc (1 + \sigma_k^{-1})^{-1}} \quad (6)$$

Proof. The proof follows from algebraic manipulations relating Equation (4) and Equation (5) from Lemma 3. ■

5 Numerical Experiments

In this section we will test various Synthetic Matrices, Differential Operators, Images, and real-world applications, with our framework compared to fixed covariance matrices. In our first experiment we attempt to learn the discretized 250×250 matrix of the inverse of the following differential operator:

$$\mathcal{L}u = \frac{\partial^2 u}{\partial x^2} - 100 \sin(5\pi x) u, \quad x \in [0, 1] \quad (7)$$

Learning the inverse operator of a PDE is equivalent to learning the Green's Function of a PDE. This has been theoretically proven for certain classes of PDEs (Linear Parabolic [5, 3]) as the inverse Differential operator is compact and there are nice theoretical properties, such as data efficiency.

In Figure 1(Right), note if the Covariance Matrix has eigenvectors orthogonal to the left singular vectors of \mathbf{A} , then the randomized SVD will not perform well. Furthermore, in Figure 1(Left), we can note even without knowledge of the Green's Function, our method achieves lower error than with the Prior Covariance. We also test our algorithm against various Sparse Matrices in the Texas A& M Sparse Matrix Suite, [8]. In Figure 2 (Left), we choose a fluid dynamics problem due to its relevance in low-rank approximation [7]. The synthetic matrix is developed in the following scheme:

$$\mathbf{A} = \sum_{i=1}^{\rho} \frac{100i^\ell}{n} \mathbf{U}_{(:,i)} \mathbf{V}_{(i,:)}^H, \quad \mathbf{U} \in \mathbb{O}_{m,k}, \mathbf{V} \in \mathbb{O}_{n,k} \quad (8)$$

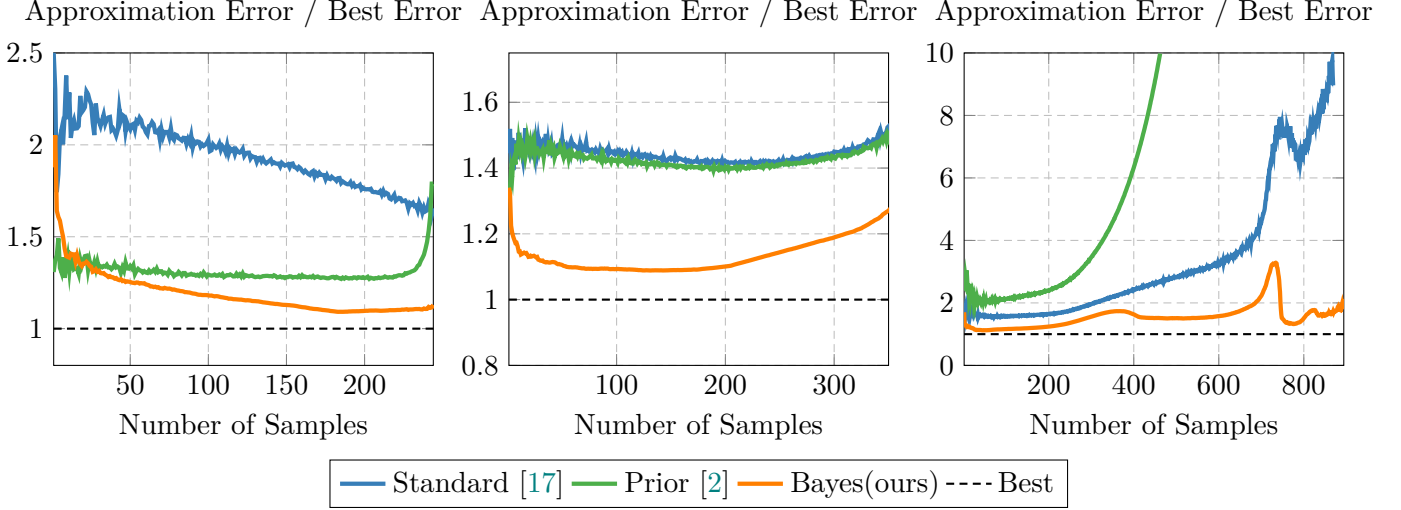


Figure 1: Low Rank Approximation for the Inverse Differential Operator given in Equation (7) (*Left*), Differential Operator Matrix `Poisson2D` [8] (*Center*), and Differential Operator Matrix `DK01R` [8] (*Right*). The experiment on the left is from [2, Figure 2].

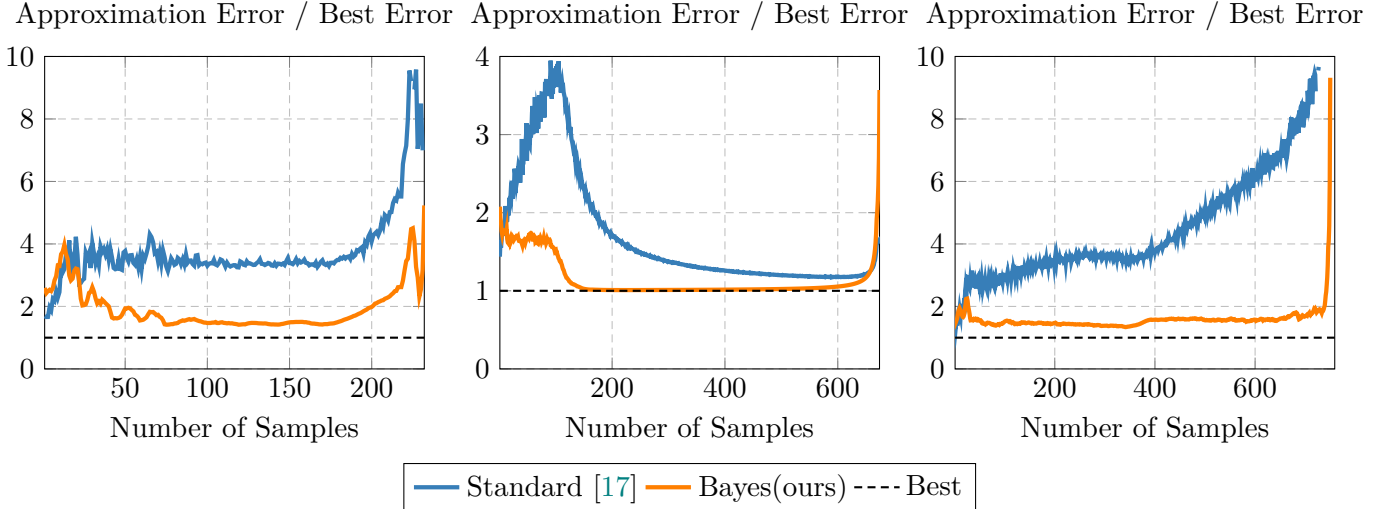


Figure 2: Low Rank Approximation for a matrix for a Computational Fluid Dynamics Problem, `saylr1` (*Left*) from [8]. Subsequent 2D/3D Problem `fs-680-2` (*Center*) from [8]. Astrophysics 2D/3D Problem `msfe` (*Right*) from [8].

6 Conclusions

We have theoretically and empirically analyzed a novel Covariance Update to iteratively construct the sampling matrix, Ω in the Randomized SVD algorithm. Our covariance update for generating sampling vectors and functions can find use various PDE learning applications, [4, 7]. Numerical Experiments indicate without prior knowledge of the matrix, we are able to obtain superior performance to the Randomized SVD and generalized Randomized SVD with covariance matrix utilizing prior information of the PDE. Theoretically, we provide an analysis of our update extended

to k -steps and show in expectation, under certain singular value decay conditions, we obtain better performance expectation.

Acknowledgments

We thank mentors Christopher Wang and Nicolas Boullé and supervisor Alex Townsend for the idea of extending Adaptive Sampling for the Matrix-Vector Product Model and the numerous helpful discussions leading to the formulation of the algorithm and the development of the theory. We also would like to thank Alex Gittens for his helpful discussions and encouragement.

References

- [1] Ainesh Bakshi, Kenneth L. Clarkson, and David P. Woodruff. Low-rank approximation with $1/\varepsilon^3$ matrix-vector products. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2022, page 11301143, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392648. doi: 10.1145/3519935.3519988. URL <https://doi.org/10.1145/3519935.3519988>.
- [2] Nicolas Boullé and Alex Townsend. A generalization of the randomized singular value decomposition. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=hgKtwSb4S2>.
- [3] Nicolas Boullé and Alex Townsend. Learning elliptic partial differential equations with randomized linear algebra. *Foundations of Computational Mathematics*, 23(2):709–739, Apr 2023. ISSN 1615-3383. doi: 10.1007/s10208-022-09556-w. URL <https://doi.org/10.1007/s10208-022-09556-w>.
- [4] Nicolas Boullé, Christopher J. Earls, and Alex Townsend. Data-driven discovery of green’s functions with human-understandable deep learning. *Scientific Reports*, 12(1):4824, Mar 2022. ISSN 2045-2322. doi: 10.1038/s41598-022-08745-5. URL <https://doi.org/10.1038/s41598-022-08745-5>.
- [5] Nicolas Boullé, Seick Kim, Tianyi Shi, and Alex Townsend. Learning green’s functions associated with time-dependent partial differential equations. *The Journal of Machine Learning Research*, 23(1):9797–9830, 2022.
- [6] Mark Braverman, Elad Hazan, Max Simchowitz, and Blake Woodworth. The gradient complexity of linear regression. In *Conference on Learning Theory*, pages 627–647. PMLR, 2020.
- [7] Steven L Brunton, Bernd R Noack, and Petros Koumoutsakos. Machine learning for fluid mechanics. *Annual review of fluid mechanics*, 52:477–508, 2020.
- [8] Timothy A. Davis and Yifan Hu. The university of florida sparse matrix collection. *ACM Trans. Math. Softw.*, 38(1), dec 2011. ISSN 0098-3500. doi: 10.1145/2049662.2049663. URL <https://doi.org/10.1145/2049662.2049663>.
- [9] Amit Deshpande and Santosh Vempala. Adaptive sampling and fast low-rank matrix approximation. In *International Workshop on Approximation Algorithms for Combinatorial Optimization*, pages 292–303. Springer, 2006.

- [10] Amit Deshpande, Luis Rademacher, Santosh S Vempala, and Grant Wang. Matrix approximation and projective clustering via volume sampling. *Theory of Computing*, 2(1):225–247, 2006.
- [11] Yijun Dong, Per-Gunnar Martinsson, and Yuji Nakatsukasa. Efficient bounds and estimates for canonical angles in randomized subspace approximations. *arXiv preprint arXiv:2211.04676*, 2022.
- [12] Petros Drineas, Ilse C. F. Ipsen, Eugenia-Maria Kontopoulou, and Malik Magdon-Ismael. Structural convergence results for approximation of dominant subspaces from block krylov spaces. *SIAM Journal on Matrix Analysis and Applications*, 39(2):567–586, 2018. doi: 10.1137/16M1091745. URL <https://doi.org/10.1137/16M1091745>.
- [13] Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- [14] Hui-Yuan Fan, George S Dulikravich, and Zhen-Xue Han. Aerodynamic data modeling using support vector machines. *Inverse Problems in Science and Engineering*, 13(3):261–278, 2005.
- [15] Alan Frieze, Ravi Kannan, and Santosh Vempala. Fast monte-carlo algorithms for finding low-rank approximations. *Journal of the ACM (JACM)*, 51(6):1025–1041, 2004.
- [16] Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- [17] N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, 2011. doi: 10.1137/090771806. URL <https://doi.org/10.1137/090771806>.
- [18] Sarel Har-Peled. Low rank matrix approximation in linear time. *arXiv preprint arXiv:1410.8802*, 2014.
- [19] Moritz Hardt and Eric Price. The noisy power method: A meta algorithm with applications. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. URL https://proceedings.neurips.cc/paper_files/paper/2014/file/729c68884bd359ade15d5f163166738a-Paper.pdf.
- [20] O. Hölder. Ueber einen mittelwerthabsatz. *Nachrichten von der Königl. Gesellschaft der Wissenschaften und der Georg-Augusts-Universität zu Göttingen*, 1889:38–47, 1889. URL <http://eudml.org/doc/180218>.
- [21] Harvard Lomax, Thomas H Pulliam, David W Zingg, Thomas H Pulliam, and David W Zingg. *Fundamentals of computational fluid dynamics*, volume 246. Springer, 2001.
- [22] Yuetian Luo, Rungang Han, and Anru R Zhang. A schatten-q low-rank matrix perturbation analysis via perturbation projection error bound. *Linear Algebra and its Applications*, 630: 225–240, 2021.
- [23] Per-Gunnar Martinsson and Joel A. Tropp. Randomized numerical linear algebra: Foundations and algorithms. *Acta Numerica*, 29:403572, 2020. doi: 10.1017/S0962492920000021.
- [24] L. Mirsky. Symmetric gauge functions and unitarily invariant norms. *The Quarterly Journal of Mathematics*, 11(1):50–59, 01 1960. ISSN 0033-5606. doi: 10.1093/qmath/11.1.50. URL <https://doi.org/10.1093/qmath/11.1.50>.

- [25] Leon Mirsky. Symmetric gauge functions and unitarily invariant norms. *The quarterly journal of mathematics*, 11(1):50–59, 1960.
- [26] Sean O’Rourke, Van Vu, and Ke Wang. Random perturbation of low rank matrices: Improving classical bounds. *Linear Algebra and its Applications*, 540:26–59, 2018. doi: 10.1016/j.laa.2017.11.014. URL <https://doi.org/10.1016/j.laa.2017.11.014>.
- [27] Saurabh Paul, Malik Magdon-Ismail, and Petros Drineas. Column selection via adaptive sampling. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/d395771085aab05244a4fb8fd91bf4ee-Paper.pdf.
- [28] Arvind K Saibaba. Randomized subspace iteration: Analysis of canonical angles and unitarily invariant norms. *SIAM Journal on Matrix Analysis and Applications*, 40(1):23–48, 2019.
- [29] Ruo-Chun Tzeng, Po-An Wang, Florian Adriaens, Aristides Gionis, and Chi-Jen Lu. Improved analysis of randomized svd for top-eigenvector approximation. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 2045–2072. PMLR, 28–30 Mar 2022. URL <https://proceedings.mlr.press/v151/tzeng22a.html>.
- [30] Roman Vershynin. High-dimensional probability. *University of California, Irvine*, 2020.
- [31] Shipeng Wang, Xiaorong Li, Jian Sun, and Zongben Xu. Training networks in null space of feature covariance for continual learning. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 184–193, 2021. doi: 10.1109/CVPR46437.2021.00025. URL <http://dx.doi.org/10.1109/CVPR46437.2021.00025>.
- [32] Per-Åke Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111, Mar 1972. ISSN 1572-9125. doi: 10.1007/BF01932678. URL <https://doi.org/10.1007/BF01932678>.
- [33] David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.

A Deferred Proofs of Main Results

In this section we give proofs for results we deferred from the main text.

A.1 Proof of Lemma 3

We have \mathbf{Q} is an orthonormal basis of $\mathbf{Y} \in \mathbb{C}^{m \times k}$ where $\mathbf{Y} \triangleq \mathbf{A}\mathbf{\Omega}$ for $\mathbf{\Omega} \in \mathbb{C}^{n \times k}$ is an arbitrary test matrix. Then let us denote $\hat{\mathbf{Q}} \triangleq \text{orth}([\mathbf{Y} \quad \mathbf{A}\mathbf{v}])$ where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\tilde{\mathbf{Q}} \triangleq \text{orth}([\mathbf{Y} \quad \mathbf{A}\tilde{\mathbf{v}}])$ where $\tilde{\mathbf{v}}$ is the k th right singular vector of $\mathbf{Q}\mathbf{Q}^H\mathbf{A}$. Note $\hat{\mathbf{q}} \in \text{Span}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H\mathbf{A})$.

$$\|\mathbf{A} - \hat{\mathbf{Q}}\hat{\mathbf{Q}}^H\mathbf{A}\|_F^2 = \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H\mathbf{A} - \hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{A}\| \quad (9)$$

$$= \text{Tr} \left(\mathbf{A}^H\mathbf{A} - \mathbf{A}^H\mathbf{Q}\mathbf{Q}^H\mathbf{A} - \mathbf{A}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{A} - \mathbf{A}^H\mathbf{Q}\mathbf{Q}^H\mathbf{A} + \mathbf{A}^H\mathbf{Q}\mathbf{Q}^H\mathbf{Q}\mathbf{Q}^H\mathbf{A} \right. \\ \left. + \mathbf{A}^H\mathbf{Q}\mathbf{Q}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{A} - \mathbf{A}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{A} + \mathbf{A}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{Q}\mathbf{Q}^H\mathbf{A} + \mathbf{A}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{A} \right) \quad (10)$$

$$= \text{Tr} \left(\mathbf{A}^H\mathbf{A} - \mathbf{A}^H\mathbf{Q}\mathbf{Q}^H\mathbf{A} - \mathbf{A}^H\hat{\mathbf{q}}\hat{\mathbf{q}}^H\mathbf{A} \right) \quad (11)$$

$$= \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^H\mathbf{A}\|_F^2 - \underbrace{\|\hat{\mathbf{q}}^H\mathbf{A}\|_F^2}_{c_2} \quad (12)$$

Similarly from Equation (12), we have

$$\|\mathbf{A} - \tilde{\mathbf{Q}}\tilde{\mathbf{Q}}^H\mathbf{A}\|_F^2 = \|\mathbf{A}\|_F^2 - \|\mathbf{Q}^H\mathbf{A}\|_F^2 - \underbrace{\|\tilde{\mathbf{q}}^H\mathbf{A}\|_F^2}_{c_1} \quad (13)$$

Since we have $\hat{\mathbf{q}}, \tilde{\mathbf{q}} \in \text{Span}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H\mathbf{A})$, then from our formulation in Equations (12) and (13), we want to our sampled vector to be in the dominant singular space of the span of the singular vectors of $\mathbf{I} - \mathbf{Q}\mathbf{Q}^H\mathbf{A}$.

$$\mathbf{q}_{\text{OPT}} = \arg \max_{\mathbf{q} \in \text{Span}(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H): \|\mathbf{q}\|=1} \|\mathbf{q}^H\mathbf{A}\| \quad (14)$$

$$= \arg \max_{\mathbf{v} \in \mathbb{R}^n} \frac{\mathbf{A}^H (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) (\mathbf{A}\mathbf{v})}{\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) (\mathbf{A}\mathbf{v})\|} \quad (15)$$

Let us note for any column $\mathbf{q} \in \mathbf{Q}$, we have

$$\left((\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{v} \right)^H \mathbf{q} = (\mathbf{v}^H - \mathbf{v}^H\mathbf{Q}\mathbf{Q}^H) \mathbf{q} = \mathbf{v}^H\mathbf{q} - \mathbf{v}^H\mathbf{q} = 0 \quad (16)$$

Claim (i):

Let us define $\tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^H = \text{SVD}(\mathbf{Q}\mathbf{Q}^H\mathbf{A})$.

$$c_1 \triangleq \|\tilde{\mathbf{q}}^H\mathbf{A}\|_F^2 \stackrel{(16)}{=} \frac{\left\| \left((\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\tilde{\mathbf{v}} \right)^H \mathbf{A} \right\|_2^2}{\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\tilde{\mathbf{v}}\|_2^2} \stackrel{(i)}{\geq} \frac{\|\tilde{\mathbf{v}}^H\mathbf{A}^H (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|_2^2}{\|\mathbf{A}\tilde{\mathbf{v}}_k\|_2^2} \quad (17)$$

$$\stackrel{\text{lem. 8}}{\geq} \frac{\|\tilde{\mathbf{v}}^H\mathbf{A}^H\mathbf{A} - \tilde{\mathbf{u}}_k^H\mathbf{A}\|_2^2}{\sigma_k^2} \stackrel{\text{lem. 8}}{\geq} \frac{\|\tilde{\mathbf{v}}_k^H\mathbf{A}^H\mathbf{A}\|_2^2 + \|\tilde{\mathbf{u}}_k^H\mathbf{A}\|_2^2 - 2\|\tilde{\mathbf{v}}_k^H\mathbf{A}^H\mathbf{A}\| \|\tilde{\mathbf{u}}_k\mathbf{A}\|}{\sigma_k^2} \quad (18)$$

$$\stackrel{\text{lem. 8, lem. 7}}{\geq} \frac{\sigma_k^4 (1 - 2 \sin^2 \Theta(\mathbf{v}_k, \tilde{\mathbf{v}}_k)) + \sigma_k^2 (1 - 2 \sin^2 \Theta(\mathbf{u}_k, \tilde{\mathbf{u}}_k)) - 2\sigma_k^3}{\sigma_k^2} \quad (19)$$

$$\stackrel{\text{thm. 9}}{\geq} (\sigma_k - 1)^2 - 2\sigma_k^2 \frac{\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A} \tilde{\mathbf{v}}_k\|_2^2}{\sigma_k^2} - 2 \frac{\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}^H \tilde{\mathbf{u}}_k\|_2^2}{\sigma_k^2} \quad (20)$$

$$\stackrel{(ii)}{\geq} (\sigma_k - 1)^2 - 2c \left(\sigma_{k+1}^2 + \frac{\sigma_{k+1}^2}{\sigma_k^2} \right) \quad (21)$$

(i) follows from noting that Π is a contraction. (ii) follows from Lemma 6 which has been seen in [12]. We thus have

$$\|\mathbf{A} - \tilde{\mathbf{Q}}\tilde{\mathbf{Q}}^H \mathbf{A}\|_F^2 \stackrel{(21)}{\leq} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|_F^2 - (\sigma_k - 1)^2 + 2c\sigma_{k+1}^2 (1 + \sigma_k^{-2}) \quad (22)$$

$$= C^2 \sigma_{k+1}^2 \left(1 - \frac{(\sigma_k - 1)^2 - 2c\sigma_{k+1}^2 (1 + \sigma_k^{-2})}{C^2 \sigma_{k+1}^2} \right) \quad (23)$$

Taking the square root of both sides completes our proof for Claim (i). ■

Claim (ii):

We will lower bound c_2 .

$$c_2 \triangleq \|\hat{\mathbf{q}}^H \mathbf{A}\|_F^2 \stackrel{(16)}{=} \frac{\|((\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A} \hat{\mathbf{v}})^H \mathbf{A}\|_2^2}{\|((\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A} \hat{\mathbf{v}})\|_2^2} \quad (24)$$

$$= \frac{\|\hat{\mathbf{v}}^H \mathbf{A}^H (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|_2^2}{\|((\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A} \hat{\mathbf{v}})\|_2^2} \geq \frac{\|\hat{\mathbf{v}}^H \mathbf{A}^H (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|_2^2}{\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|_2^2 \|\hat{\mathbf{v}}\|_2^2} \quad (25)$$

For shorthand, let $\tilde{\mathbf{A}} \triangleq \mathbf{A}^H (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}$. Note, we have $\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|$ is given. From Lemma 12, we have

$$\mathbb{P} \left\{ \left| \mathbf{x}^H \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \mathbf{x} - \mathbb{E} \left[\mathbf{x}^H \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \mathbf{x} \right] \right| \geq t \right\} \leq 2 \exp \left(-c \frac{t^2}{\|\tilde{\mathbf{A}}^H \tilde{\mathbf{A}}\|_F^2} \right) \quad (26)$$

It then follows with probability $1 - \delta$

$$\|\tilde{\mathbf{A}} \mathbf{x}\|^2 \geq \mathbb{E} \|\tilde{\mathbf{A}} \mathbf{x}\|^2 - \|\tilde{\mathbf{A}}\|_F^2 \sqrt{\frac{1}{c} \log \left(\frac{2}{\delta} \right)} \quad (27)$$

We further simplify this with the fact

$$\mathbb{E} \|\tilde{\mathbf{A}} \mathbf{x}\|^2 = \mathbb{E} \mathbf{x}^H \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \mathbf{x} = \mathbb{E} \text{Tr} \left(\mathbf{x}^H \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \mathbf{x} \right) = \mathbb{E} \text{Tr} \left(\mathbf{x} \mathbf{x}^H \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \right) = \text{Tr} \left(\mathbb{E} \left[\mathbf{x} \mathbf{x}^H \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \right] \right) \quad (28)$$

$$= \text{Tr} \left(\mathbb{E} \left[\mathbf{x} \mathbf{x}^H \right] \tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \right) = \text{Tr} \left(\tilde{\mathbf{A}}^H \tilde{\mathbf{A}} \right) = \|\tilde{\mathbf{A}}\|_F^2 \quad (29)$$

Let us also note the following fact for a rank- k orthonormal matrix, \mathbf{Q} , and rank- r matrix \mathbf{A} .

$$\|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}\|_F^2 = \text{Tr} \left(\mathbf{A} (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}^H \right) \quad (30)$$

$$= \text{Tr} \left(\mathbf{A} (\mathbf{I} - \mathbf{Q}\mathbf{Q}^H) \mathbf{A}^H \right) \quad (31)$$

$$\stackrel{\zeta_1}{\leq} \sqrt{r} \left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A}^H \right\|_F \quad (32)$$

where (ζ_1) follows from the relation of the $\|\cdot\|_1$ and $\|\cdot\|_2$ norm. Now we can calculate the expectation. We then note $\mathbb{E} \|\mathbf{v}\|_2^2$ is the expectation of n -degree of freedom Chi-squared variable and thus is equal to n , then

$$\begin{aligned} \mathbb{E} c_2 &= \mathbb{E} \left[\frac{\left\| \hat{\mathbf{v}}^H \mathbf{A}^H \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_2^2}{\left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \hat{\mathbf{v}} \right\|_2^2} \right] \stackrel{(\zeta_1)}{\geq} \left(\mathbb{E} \left\| \hat{\mathbf{v}}^H \mathbf{A}^H \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_2 \right)^2 \mathbb{E} \left[\left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \hat{\mathbf{v}} \right\|_2^2 \right]^{-1} \\ &\stackrel{(\zeta_2)}{\geq} \left(\frac{16}{75\sqrt{5}} \right)^2 \left\| \mathbf{A}^H \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^2 \left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^{-2} \end{aligned} \quad (33)$$

$$\stackrel{(30)}{\geq} \left(\frac{16}{75\sqrt{5}} \right)^2 \frac{\left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^4}{r \left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^2} \gtrsim \left(\frac{1}{r} \right) \left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^2 \quad (34)$$

(ζ_1) follows from the Reverse Hölder Inequality [20]. (ζ_2) follows from an application of Jensen's Inequality. Combining Equation (34) and Equation (29) we thus have with probability $1 - \delta$

$$c_2 \geq \left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^2 \left(\Omega \left(\frac{1}{r} \right) \left(1 - \sqrt{\frac{1}{c_3} \log \left(\frac{2}{\delta} \right)} \right) \right) \quad (35)$$

Furthermore, in expectation we have

$$\mathbb{E} c_2 \geq \Omega \left(\frac{1}{r} \right) \left\| \left(\mathbf{I} - \mathbf{Q}\mathbf{Q}^H \right) \mathbf{A} \right\|_F^2 \quad (36)$$

This completes our proof of Claim (ii). ■

B Singular Subspace Perturbation Lemmas

Lemma 5. *Let \mathbf{v} and $\tilde{\mathbf{v}}$ be vectors s.t. $\|\mathbf{v}\| = \|\tilde{\mathbf{v}}\| = 1$ and $\mathbf{v}^H \tilde{\mathbf{v}} \geq 0$. Then,*

$$\|\mathbf{v} - \tilde{\mathbf{v}}\| \leq \sqrt{2} \sin \Theta(\mathbf{v}, \tilde{\mathbf{v}}) \quad (37)$$

Proof.

$$\sin^2 \Theta(\mathbf{v}, \tilde{\mathbf{v}}) = 1 - \left(\mathbf{v}^H \tilde{\mathbf{v}} \right)^2 \stackrel{(a)}{\geq} 1 - \mathbf{v}^H \tilde{\mathbf{v}} = 1 + \frac{1}{2} \|\mathbf{v} - \tilde{\mathbf{v}}\|^2 - \frac{1}{2} \|\mathbf{v}\|^2 - \frac{1}{2} \|\tilde{\mathbf{v}}\|^2 = \frac{1}{2} \|\mathbf{v} - \tilde{\mathbf{v}}\|^2 \quad (38)$$

(a) follows from $0 \leq \mathbf{v}^H \tilde{\mathbf{v}} \leq 1$, therefore $\mathbf{v}^H \tilde{\mathbf{v}} \geq \left(\mathbf{v}^H \tilde{\mathbf{v}} \right)^2$.

Plugging this back into the first inequality and taking the square root gives us the desired result. ■

Lemma 6. *Let $\tilde{\mathbf{u}}_k$ represent the k th left singular vector of a low-rank approximation $\tilde{\mathbf{A}}$ of \mathbf{A} . Then,*

$$\left\| \mathbf{A}^H \tilde{\mathbf{u}}_k \right\| \leq \sigma_k(\mathbf{A}) \quad (39)$$

Proof. We provide the proof for completeness. From the minimax Courant-Fischer Theorem [16] we have,

$$\sigma_k(\mathbf{A}) = \max_{\dim(S)=k} \min_{\mathbf{x} \in S \setminus \{\mathbf{0}\}} \frac{\left\| \mathbf{A}^H \mathbf{x} \right\|}{\|\mathbf{x}\|} \quad (40)$$

$$\stackrel{\zeta_1}{=} \min_{\mathbf{x} \in \text{Span}(\{\mathbf{u}_1, \dots, \mathbf{u}_k\}) \setminus \{\mathbf{0}\}} \frac{\|\mathbf{A}^H \mathbf{x}\|}{\|\mathbf{x}\|} \quad (41)$$

$$\geq \min_{\mathbf{x} \in \text{Span}(\{\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_k\}) \setminus \{\mathbf{0}\}} \frac{\|\mathbf{A}^H \mathbf{x}\|}{\|\mathbf{x}\|} \quad (42)$$

$$= \|\mathbf{A}^H \tilde{\mathbf{u}}_k\| \quad (43)$$

where (ζ_1) follows from the optimality of the SVD given by the Eckart-Young-Mirsky Theorem [13, 25].

For $k = 1$, we have

$$\|\mathbf{A}^H \tilde{\mathbf{u}}_1\| = \|\mathbf{V} \Sigma \mathbf{U}^H \mathbf{u}_1\| = \|\Sigma \mathbf{U}^H \mathbf{u}_1\| \leq \sigma_1 \cos(\angle \mathbf{u}_1, \mathbf{u}_1) + \sigma_2 \sin \Theta(\mathbf{u}_1, \tilde{\mathbf{u}}_1) \quad (44)$$

■

Lemma 7. *Let $\tilde{\mathbf{v}}_k$ be the k th right singular vector of an approximation $\tilde{\mathbf{A}}$ of \mathbf{A} . Then,*

$$\|\mathbf{A}^H \mathbf{A} \tilde{\mathbf{v}}_k\| \geq \sigma_k^2 (1 - \sqrt{2} \sin \Theta(\mathbf{v}_k, \tilde{\mathbf{v}}_k)) \quad (45)$$

Proof.

$$\|\mathbf{A}^H \mathbf{A}\|_{\text{F}}^2 - \|\mathbf{A}^H \mathbf{A} \tilde{\mathbf{v}}_k\|_{\text{F}}^2 = \|\mathbf{A}^H \mathbf{A} - \tilde{\mathbf{v}}_k \tilde{\mathbf{v}}_k^H \mathbf{A}^H \mathbf{A}\|_{\text{F}}^2 \stackrel{\zeta_1}{=} \|\mathbf{A}_{(k)}^H \mathbf{A}_{(k)} - \tilde{\mathbf{v}}_k \tilde{\mathbf{v}}_k^H \mathbf{A}^H \mathbf{A}\|_{\text{F}}^2 + \|\mathbf{A}_{\perp, k}^H \mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (46)$$

$$\leq \|\mathbf{A}_{(k)}^H \mathbf{A}_{(k)} - \sigma_k^2 \tilde{\mathbf{v}}_k \tilde{\mathbf{v}}_k^H\|_{\text{F}}^2 + \|\mathbf{A}_{\perp, k}^H \mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (47)$$

$$= \sigma_k^4 \|\Pi_{\mathbf{v}_k} - \Pi_{\tilde{\mathbf{v}}_k}\|_{\text{F}}^2 + \|\mathbf{A}_{\perp, k}^H \mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (48)$$

$$\stackrel{\text{lem. 5}}{=} 2\sigma_k^4 \sin^2 \Theta(\mathbf{v}_k, \tilde{\mathbf{v}}_k) + \|\mathbf{A}_{\perp, k}^H \mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (49)$$

(ζ_1) follows from the Matrix Pythagoras theorem [33]. From rearranging the inequalities and noting

$\|\mathbf{A}^H \mathbf{A}\|_{\text{F}}^2 - \|\mathbf{A}_{\perp, k}^H \mathbf{A}_{\perp, k}\|_{\text{F}}^2 = \sigma_k^4$, we then obtain

$$\|\mathbf{A}^H \mathbf{A} \tilde{\mathbf{v}}_k\|_{\text{F}}^2 \geq \sigma_k^4 (1 - 2 \sin^2 \Theta(\mathbf{v}_k, \tilde{\mathbf{v}}_k)) \quad (50)$$

Taking the square root and reverse triangle inequality gives us the desired result. ■

Lemma 8. *Let $\tilde{\mathbf{u}}_k$ be the k th left singular vector of approximation $\tilde{\mathbf{A}}$ of \mathbf{A} , then we have*

$$\|\mathbf{A}^H \tilde{\mathbf{u}}_k\| \geq \sigma_k (1 - \sqrt{2} \sin \Theta(\mathbf{u}_k, \tilde{\mathbf{u}}_k)) \quad (51)$$

Proof.

$$\|\mathbf{A}\|_{\text{F}}^2 - \|\tilde{\mathbf{u}}_k^H \mathbf{A}\|_{\text{F}}^2 = \|\mathbf{A} - \tilde{\mathbf{u}}_k \tilde{\mathbf{u}}_k^H \mathbf{A}\|_{\text{F}}^2 \quad (52)$$

$$\leq \|\mathbf{A} - \tilde{\mathbf{u}}_k \tilde{\mathbf{u}}_k^H \mathbf{A}_{(k)}\|_{\text{F}}^2 \quad (53)$$

$$= \|\mathbf{A}_{(k)} - \tilde{\mathbf{u}}_k \tilde{\mathbf{u}}_k^H \mathbf{A}_{(k)}\|_{\text{F}}^2 + \|\mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (54)$$

$$\leq \|\mathbf{A}_{(k)} - \sigma_k \tilde{\mathbf{u}}_k \mathbf{v}_k^H\|_{\text{F}}^2 + \|\mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (55)$$

$$\leq 2\sigma_k^2 \sin^2 \Theta(\mathbf{u}_k, \tilde{\mathbf{u}}_k) + \|\mathbf{A}_{\perp, k}\|_{\text{F}}^2 \quad (56)$$

Rearranging we obtain the desired inequality. ■

Theorem 9. [32]. Let $\mathbf{A}, \hat{\mathbf{A}} \in \mathbb{R}^{m \times n}$ have singular values $\sigma_1 \geq \dots \geq \sigma_{m \vee n}$ and $\hat{\sigma}_1 \geq \dots \geq \hat{\sigma}_{m \vee n}$, respectively. Given $j \in 1, \dots, m \vee n$,

$$\sin \Theta(\mathbf{v}_1, \tilde{\mathbf{v}}_1) \leq \frac{\|\mathbf{A} - \hat{\mathbf{A}}\|}{\sigma_1 - \hat{\sigma}_2} \quad (57)$$

Theorem 10. [26]. Let $\mathbf{A}, \hat{\mathbf{A}} \in \mathbb{R}^{m \times n}$ have singular values $\sigma_1 \geq \dots \geq \sigma_{m \vee n}$ and $\hat{\sigma}_1 \geq \dots \geq \hat{\sigma}_{m \vee n}$, respectively. Then,

$$\sin \Theta(\mathbf{v}_1, \tilde{\mathbf{v}}_1) \leq \frac{2 \|\mathbf{A} - \hat{\mathbf{A}}\|}{\sigma_1 - \sigma_2} \quad (58)$$

C Concentration Inequalities

Lemma 11 ([30], Exercise 6.3.5). Let $\mathbf{B} \in \mathbb{C}^{m \times n}$ and $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then we have

$$\mathbb{P}\{\|\mathbf{B}\mathbf{x}\|_2 \geq CK \|\mathbf{B}\|_F + t\} \leq \exp\left(-\frac{ct^2}{K \|\mathbf{B}\|^2}\right) \quad (59)$$

Lemma 12 (Hanson-Wright Inequality). Let $\mathbf{x} \in \mathbb{R}^n$ be a vector with sub-Gaussian random vector symmetric about $\mathbf{0}$. Let \mathbf{B} be a symmetric $n \times n$ matrix, then $\forall t \geq 0$,

$$\mathbb{P}\left\{\left|\mathbf{x}^\top \mathbf{B} \mathbf{x} - \mathbb{E}[\mathbf{x}^\top \mathbf{B} \mathbf{x}]\right| \geq t\right\} \leq 2 \exp\left(-c \min\left\{\frac{t^2}{K^4 \|\mathbf{B}\|_F^2}, \frac{t}{K^2 \|\mathbf{B}\|}\right\}\right) \quad (60)$$

Lemma 13. [24]. For any matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} , then for any unitarily invariant norm $\|\cdot\|$, it follows

$$\|\mathbf{A}\mathbf{B}\mathbf{C}\| \leq \min\{\|\mathbf{A}\| \|\mathbf{B}\|_2 \|\mathbf{C}\|_2, \|\mathbf{A}\|_2 \|\mathbf{B}\| \|\mathbf{C}\|_2, \|\mathbf{A}\|_2 \|\mathbf{B}\|_2 \|\mathbf{C}\|\} \quad (61)$$

Lemma 14. Let $\mathbf{A} \in \mathbb{C}^{m \times n}$ and $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, then we have

$$\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2 \geq \frac{16}{75\sqrt{5}} \|\mathbf{A}\|_F \quad (62)$$

Proof. For a $\theta \in (0, 1)$, we have

$$\frac{\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2}{\theta \sqrt{\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2}} \stackrel{(\zeta_1)}{\geq} \mathbb{P}\left\{\|\mathbf{A}\mathbf{x}\|_2 \geq \theta \sqrt{\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2}\right\} = \mathbb{P}\left\{\|\mathbf{A}\mathbf{x}\|_2^2 \geq \theta \mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2\right\} \stackrel{(\zeta_2)}{\geq} \frac{(1 - \theta^2)^2 (\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2)^2}{\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^4} \quad (63)$$

(ζ_1) follows from Markov's Inequality. (ζ_2) follows from the Paley-Zygmund Inequality. Rearranging the LHS of Equation (63) with RHS of Equation (63), we have

$$\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2 \geq \theta(1 - \theta^2)^2 \frac{(\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2)^{5/2}}{\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^4} \stackrel{(\zeta_3)}{\geq} \frac{16 (\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2)^{5/2}}{25\sqrt{5} \mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^4} \quad (64)$$

where (ζ_3) follows from noting $\theta(1 - \theta^2)^2$ is maximized at $\theta = \frac{1}{\sqrt{5}}$. Next we note $\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{A}\|_F^2$. Furthermore, we have

$$\mathbb{E} \|\mathbf{A}\mathbf{x}\|_2^4 = \mathbb{E} \left\| \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} \right\|_2^2 = \text{Tr}(\mathbf{A}^\top \mathbf{A})^2 + 2 \text{Tr}((\mathbf{A}^\top \mathbf{A})^2) \leq \|\mathbf{A}\|_F^4 + 2 \|\mathbf{A}\|_F^2 \|\mathbf{A}\|_2^2 \leq 3 \|\mathbf{A}\|_F^4 \quad (65)$$

Then we can substitute Equation (65) into Equation (64), and we have

$$\mathbb{E} \|\mathbf{Ax}\|_2 \stackrel{(65)}{\geq} \frac{16 \left(\|\mathbf{A}\|_F^2 \right)^{5/2}}{75\sqrt{5} \|\mathbf{A}\|^4} = \frac{16}{75\sqrt{5}} \|\mathbf{A}\|_F \quad (66)$$

This concludes the proof. ■

D Additional Experiments

In this section we perform more experiments on learning the inverse operator for PDE matrices with State of the Art Matrix Experiments.

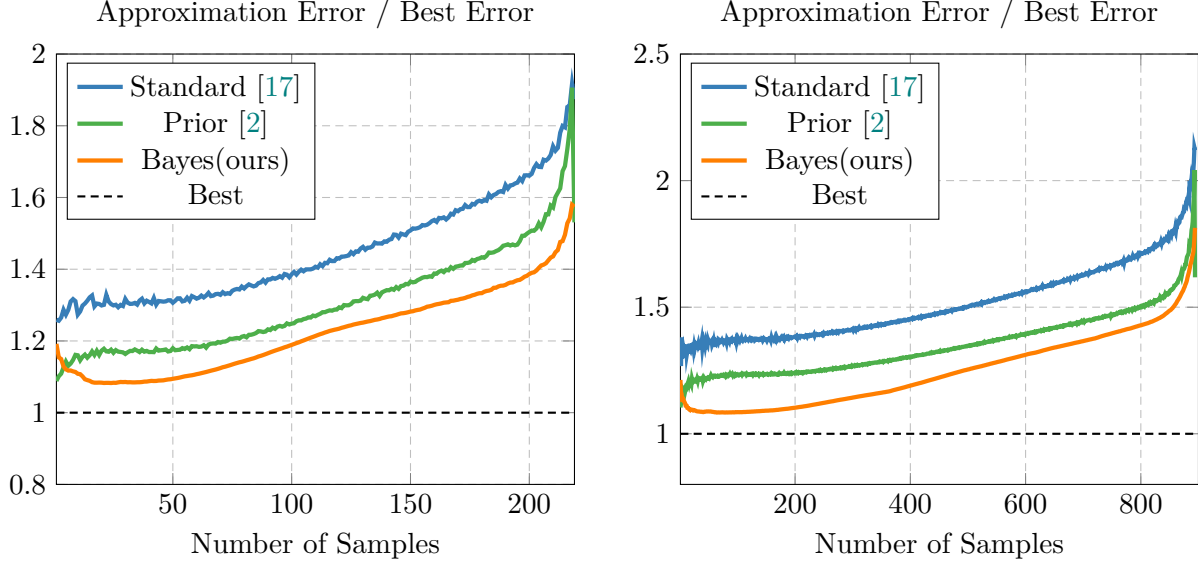


Figure 3: In (*Left*), Matrix from TAMU Sparse Matrix Suite `pde 225`. In (*Right*), Matrix from TAMU Sparse Matrix Suite `pde 900`. With the prior, we use the covariance matrix associated with the discrete Green’s Function for the Laplacian as in Equation (7).

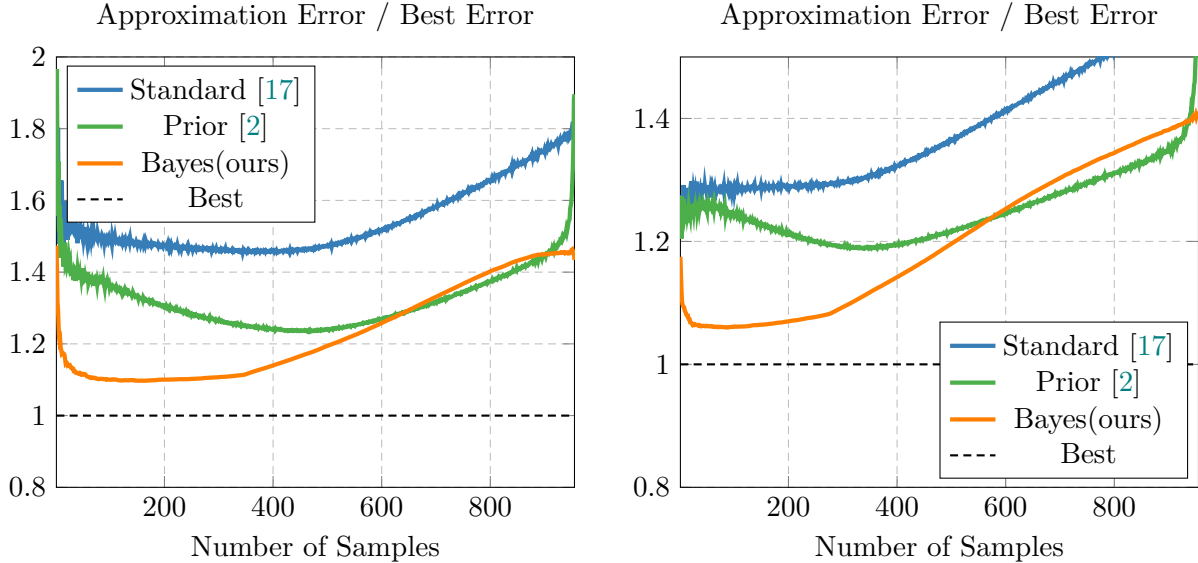


Figure 4: In these figures we look at matrices from Computational Fluid Dynamics. In (*Left*), Matrix from TAMU Sparse Matrix Suite `cdde1`. In (*Right*), Matrix from TAMU Sparse Matrix Suite `cdde1`. With the prior, we use the covariance matrix associated with the discrete Green’s Function for the Laplacian as in Equation (7).

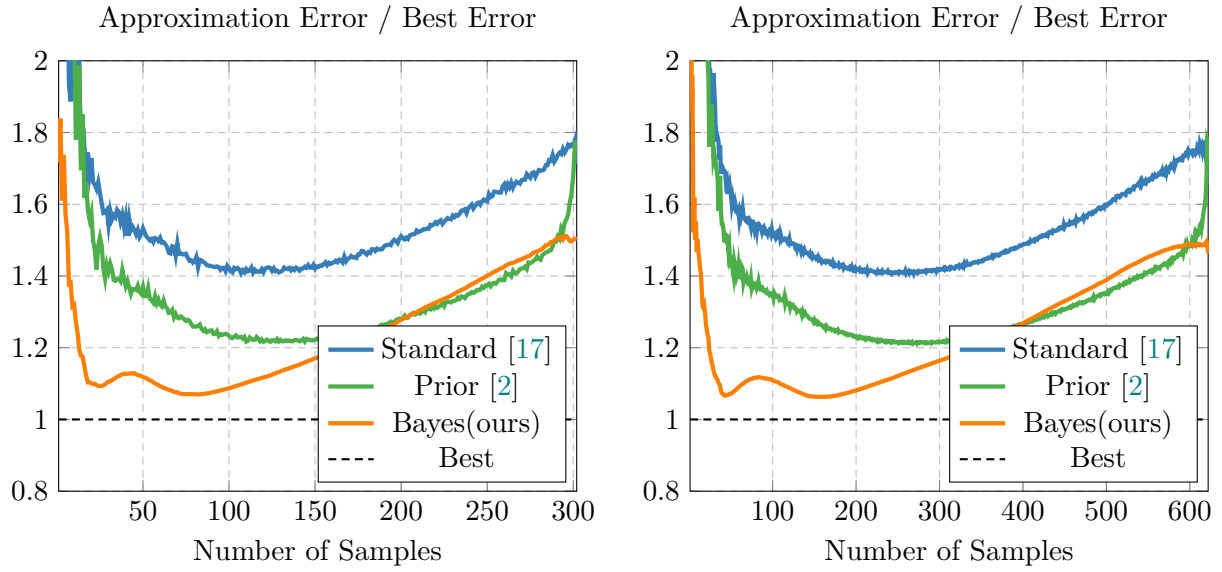


Figure 5: In these figures we look at matrices from the PDE for the Poisson Differential Operator. In (*Left*), Matrix from TAMU Sparse Matrix Suite `cz308`. In (*Right*), Matrix from TAMU Sparse Matrix Suite `cz628`. With the prior, we use the covariance matrix associated with the discrete Green's Function for the Laplacian as in Equation (7).