

STRUCTURED RANDOM SKETCHING FOR PDE INVERSE PROBLEMS

KE CHEN, QIN LI, KIT NEWTON, AND STEPHEN J. WRIGHT

ABSTRACT. For an overdetermined system $Ax \approx b$ with A and b given, the least-square (LS) formulation $\min_x \|Ax - b\|_2$ is often used to find an acceptable solution x . The cost of solving this problem depends on the dimensions of A , which are large in many practical instances. This cost can be reduced by the use of random sketching, in which we choose a matrix S with many fewer rows than A and b , and solve the sketched LS problem $\min_x \|S(Ax - b)\|_2$ to obtain an approximate solution to the original LS problem. Significant theoretical and practical progress has been made in the last decade in designing the appropriate structure and distribution for the sketching matrix S . When A and b arise from discretizations of a PDE-based inverse problem, tensor structure is often present in A and b . For reasons of practical efficiency, S should be designed to have a structure consistent with that of A . Can we claim similar approximation properties for the solution of the sketched LS problem with structured S as for fully-random S ? We give estimates that relate the quality of the solution of the sketched LS problem to the size of the structured sketching matrices, for two different structures. Our results are among the first known for random sketching matrices whose structure is suitable for use in PDE inverse problems.

1. INTRODUCTION

In overdetermined linear systems (in which the number of linear conditions exceeds the number of unknowns), the least-squares (LS) solution is often used as an approximation to the true solution when the data contains noise. Given the system $Ax = b$ where $A \in \mathbb{R}^{n \times p}$ with $n \gg p$, the least-squares solution x^* is obtained by minimizing the l^2 -norm discrepancy between the Ax and b , that is,

$$\min_x \|Ax - b\|_2, \implies x^* = A^\dagger b, \text{ where } A^\dagger \stackrel{\text{def}}{=} (A^\top A)^{-1} A^\top. \quad (1)$$

The matrix A^\dagger is often called the *pseudoinverse* (more specifically the *Moore-Penrose pseudoinverse*) of A .

The LS method is ubiquitous in statistics and engineering, but large problems can be expensive to solve. Aside from the cost of preparing A , the cost of solving for x^* is $\mathcal{O}(np^2)$ flops for general (dense) A is prohibitive in large dimensions.

We can replace the LS problem with a smaller approximate LS problem by using *sketching*. Each row of the sketched system is a linear combination of the rows of A , together with the same linear combination of the elements of b . This scheme amounts to defining a sketching matrix $S \in \mathbb{R}^{r \times n}$ with $r \ll n$, and replacing the original LS problem by

$$\min_x \|SAx - Sb\|_2, \implies x_s^* = (SA)^\dagger Sb. \quad (2)$$

For appropriate choices of S , the solutions of (1) and (2) are related in the sense that

$$\|b - Ax^*\| \text{ is not too much smaller than } \|b - Ax_s^*\|. \quad (3)$$

Usually one does not design S directly, but rather draws its entries from a certain distribution. In such a setup, we can ask whether (3) holds with high probability.

The literature on random sketching is rich. During the past decade, many theoretical and numerical studies have appeared [2, 9, 11, 13, 14, 16, 18, 19, 24–26, 31, 33, 35, 40], with applications in such subjects as stochastic optimization [18], l^p regression [10, 11, 24, 29, 31, 34, 39], and tensor decomposition [3, 4, 8, 21, 30]. The technical support for these results comes mostly from the Johnson-Lindenstrauss lemma [17], random matrix theory [36,

37], and compressed sensing [15]. Two important perspectives have been utilized. One approach starts with the least squares problem and proposes two conditions for the random matrix such that an accurate solution can be attained with high confidence. It is then shown that certain choices of random matrices indeed satisfy these two conditions. Instances of this approach can be found in [14, 31, 33] and the reviews [20, 23]. The second perspective focuses on the structure of the space spanned by \mathbf{A} . It is argued that this space can be approximated by a finite number of vectors (the so-called γ -net), which can further be “embedded” using random matrices, with high accuracy; see [10, 34, 39] and a review [40]. We use this second perspective in this paper.

There are many variations of the original sketching problem. With some statistical assumptions on the perturbation in the right hand side, results could be further enhanced [29], and the sketching problem is also investigated when other constraints (such as l_1 constraints) are present; see for example [27]. In [9, 14, 28] the authors also directly quantify $\|\mathbf{x}_s^* - \mathbf{x}^*\|$ instead of the residual, as in (3).

In most previous studies, the design of \mathbf{S} varies according to the priorities of the application. For good accuracy with small r , random projections with sub-Gaussian variables are typically used. When the priority is to reduce the cost of computing the product $\mathbf{S}\mathbf{A}$, either sparse or Hadamard type matrices have been proposed, leading to “random-sampling” or FFT-type reduction in cost of the matrix-matrix multiplication. To cure “bias” in the selection process, leverage scores have been introduced; these trace their origin back to classical methods in experimental design.

In this paper, with practical inverse problems in mind, we consider the case in which \mathbf{A} and \mathbf{b} have certain tensor-type structures. For the sketched system to be formed and solved efficiently, the random sketching matrix \mathbf{S} must have a corresponding tensor structure. For these tensor-structured sketching matrices \mathbf{S} , we ask: What are the requirements on r to achieve a certain accuracy in the solution \mathbf{x}_s^* of the sketched system?

We consider \mathbf{A} with the following structure:

$$\mathbf{A} = \mathbf{F} * \mathbf{G}, \quad (4)$$

where $*$ denotes the (column-wise) *Khatri-Rao product* of the matrices \mathbf{F} and \mathbf{G} . Assuming $i_1 \in \mathcal{I}_1$ and $i_2 \in \mathcal{I}_2$, with cardinalities $n_1 = |\mathcal{I}_1|$ and $n_2 = |\mathcal{I}_2|$, respectively, the dimensions of these matrices are

$$\mathbf{F} \in \mathbb{R}^{n_1 \times p}, \quad \mathbf{G} \in \mathbb{R}^{n_2 \times p}, \quad \mathbf{A} \in \mathbb{R}^{n \times p}, \quad (5)$$

where $n = |\mathcal{I}_1 \otimes \mathcal{I}_2| = n_1 n_2$.

By defining $\mathbf{f}_j = \mathbf{F}_{:,j} \in \mathbb{R}^{n_1}$ and $\mathbf{g}_j = \mathbf{G}_{:,j} \in \mathbb{R}^{n_2}$, we can define \mathbf{A} alternatively as

$$\mathbf{a}_j \stackrel{\text{def}}{=} \mathbf{A}_{:,j} = \mathbf{f}_j \otimes \mathbf{g}_j, \quad (6)$$

where $\mathbf{a}_j \in \mathbb{R}^n$ denotes the j th column of \mathbf{A} , for $j = 1, 2, \dots, p$. For vector \mathbf{b} , we assume that it admits the same tensor structure, that is,

$$\mathbf{b} = \mathbf{f}_b \otimes \mathbf{g}_b, \quad \text{for some fixed } \mathbf{f}_b \in \mathbb{R}^{n_1} \text{ and } \mathbf{g}_b \in \mathbb{R}^{n_2}. \quad (7)$$

This type of structure comes from the fact that to formulate inverse problems, one typically needs to prepare both the *forward* and *adjoint* solutions. Denoting by $\sigma(x)$ the unknown function to be reconstructed in the inverse PDE problem, a very typical formulation is written as a Fredholm integral of the first type:

$$\int f_{i_1}(x) g_{i_2}(x) \sigma(x) dx = \text{data}_{i_1, i_2}, \quad (8)$$

where f_{i_1} and g_{i_2} solve the forward and adjoint equations respectively, equipped with boundary/initial conditions indexed by i_1 and i_2 . Each term on the right-hand side of (8) is typically data measured at i_2 with input source index i_1 . To reconstruct σ , one loops over the entire list of conditions for f_{i_1} ($i_1 \in \mathcal{I}_1$) and g_{i_2} ($i_2 \in \mathcal{I}_2$). The LS formulation $\min \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ is the discrete version of the Fredholm integral (8).

This structure imposes requirements on the sketching matrix \mathbf{S} . Since \mathcal{I}_1 and \mathcal{I}_2 contain conditions for different sets of equations, sketching needs to be performed within \mathcal{I}_1 and \mathcal{I}_2 separately. This condition is

reflected by choosing the sketching matrix \mathbf{S} to be the *row-wise Khatri-Rao product* of \mathbf{P} and \mathbf{Q} , that is,

$$\mathbf{S}_{i,:} = \mathbf{p}_i^\top \otimes \mathbf{q}_i^\top,$$

where $\mathbf{p}_i \in \mathbb{R}^{n_1}$ and $\mathbf{q}_i \in \mathbb{R}^{n_2}$, $i = 1, \dots, p$. The product \mathbf{SA} then has the special form:

$$(\mathbf{SA})_{i,:} = (\mathbf{p}_i^\top \mathbf{F}) \circ (\mathbf{q}_i^\top \mathbf{G}), \quad \text{or equivalently} \quad (\mathbf{SA})_{i,j} = (\mathbf{p}_i^\top \mathbf{f}_j)(\mathbf{q}_i^\top \mathbf{g}_j). \quad (9)$$

Thus, to formulate the i row in the reduced (sketched) system, we perform a linear combination of parameters in \mathcal{I}_1 according to \mathbf{p}_i to feed in the forward solver, and a linear combination of parameters in \mathcal{I}_2 according to \mathbf{q}_i to feed in the adjoint solver, then assemble the results in the Fredholm integral (8).

With the structural requirements for \mathbf{S} in mind, we consider the following two approaches for choosing \mathbf{S} .

Case 1: Generate two random matrices \mathbf{P} and \mathbf{Q} , of size $r_1 \times n_1$ and $r_2 \times n_2$, respectively, and define \mathbf{S} to be their tensor product:

$$\mathbf{S} = \mathbf{P} \otimes \mathbf{Q} \in \mathbb{R}^{r_1 r_2 \times n_1 n_2}. \quad (10)$$

Case 2: Generate two sets of r random vectors $\{\mathbf{p}_i, i = 1, 2, \dots, r\}$ and $\{\mathbf{q}_i, i = 1, 2, \dots, r\}$, with $\mathbf{p}_i \in \mathbb{R}^{n_1}$ and $\mathbf{q}_i \in \mathbb{R}^{n_2}$ for each i , and define row i of \mathbf{S} to be the tensor product of the vectors \mathbf{p}_i and \mathbf{q}_i :

$$\mathbf{S} = \frac{1}{\sqrt{r}} \begin{bmatrix} \mathbf{p}_1^\top \otimes \mathbf{q}_1^\top \\ \vdots \\ \mathbf{p}_r^\top \otimes \mathbf{q}_r^\top \end{bmatrix} \in \mathbb{R}^{r \times n_1 n_2}. \quad (11)$$

Case 2 gives greater randomness, in a sense, because the rows of \mathbf{P} and \mathbf{Q} are not “re-used” as in the first option.

We are not interested in designing sketching matrices of Hadamard type. In practice, \mathbf{A} is often semi-infinite: \mathbf{F} and \mathbf{G} contain all possible forward and adjoint solutions, a set of infinite cardinality that cannot be prepared in advance. In practice, one can only obtain the “realizations” $\mathbf{p}^\top \mathbf{F}$ or $\mathbf{q}^\top \mathbf{G}$ obtained by solving the forward and adjoint equations with the parameters contained in \mathbf{p} and \mathbf{q} . Because we use this technique to find \mathbf{SA} , rather than computing the matrix-matrix product explicitly, there is no advantage to defining \mathbf{S} in terms of Hadamard type random matrices.

There have been discussions in the sketching literature on problems that share our setups, including sketching of matrices \mathbf{A} with Khatri-Rao product structure. The paper [4] presents a tensor interpolative decomposition problem which discusses Khatri-Rao product form, but there is not a focus on sketching. The paper [35] proposes a so-called tensor random projection (TRP), similar to our Case 2 presented below. However, they mainly obtain sketching of one arbitrarily given vector in the space, while we need to sketch the entire space. Directly employing their argument in our setting would lead to $r = \mathcal{O}(p^8/\varepsilon^2)$, whereas our argument suggests that having $r = \mathcal{O}(p^6/\varepsilon)$ is sufficient. This point will be discussed further in Theorem 4.

In [16, 22] the authors considered the fast Johnson-Lindenstrauss Transform (JLT) random matrices and showed that the Kronecker product of fast JLT is also a JLT. This structure allows embedding of an arbitrarily given vector. For embedding vectors that have tensor structure, [12, 13] developed *TensorSketch* or *CountSketch*, and discussed the efficiency of these algorithms in terms of the number of nonzero entries in \mathbf{A} . All these results are highly related to ours, but they all have dependences on the ambient space dimension n , making them poorly suited to our setting, where we consider the possibility of $n \rightarrow \infty$.

The rest of the paper is organized as follows. In Section 2, we give two examples from PDE-based inverse problem that give rise to a linear system with tensor structure. Section 3 presents classical results on sketching for general linear regression, and states our main results on sketching of inverse problem associated with a tensor structure. Sections 4 and 5 study the two different sketching strategies outline above. Computational testing described in Section 6 validates our results.

We denote the range space (column space) of a matrix \mathbf{X} by $\text{Range}(\mathbf{X})$.

2. OVERDETERMINED SYSTEMS WITH TENSOR STRUCTURE ARISING FROM PDE INVERSE PROBLEMS

Most PDE-based inverse problems, upon linearization, reduce to a tensor structured Fredholm integral (8), which can be discretized to formulate a sketching problem.

One particularly famous example is Electrical Impedance Tomography (EIT), in which we apply voltage strength and measure current density at the boundary of some bio-tissues to infer for conductivity inside the body. The underlying PDE is a standard second order elliptic equation

$$\begin{aligned}\nabla_x \cdot (\bar{\sigma}(x) \nabla_x \bar{\rho}(x)) &= 0, & x \in \Omega, \\ \rho(x) &= \phi(x), & x \in \partial\Omega,\end{aligned}\tag{12}$$

where $\phi(x)$ is the voltage strength applied on the surface of some bio-tissue, while $\bar{\rho}(x)$, the solution to the PDE, is the voltage generated throughout the body. The unknown conductivity $\bar{\sigma}(x)$ will be inferred. The measurements are taken on the boundary too. In particular, one measures the current density on the surface of Ω tested on a testing function ψ , as follows:

$$\overline{\text{data}}_{\phi, \psi} = \int_{\partial\Omega} \bar{\sigma}(x) \frac{\partial \bar{\rho}(x)}{\partial n} \psi(x) dx.\tag{13}$$

Here, $\frac{\partial}{\partial n}$ is the normal derivative, with n being the normal direction pointing out of domain Ω . The data has two subscripts: $\phi(x)$ is the voltage applied to the surface and $\psi(x)$ is a testing function that encodes the way measurements are taken. When the detector is extremely precise, one can set $\psi(x) = \delta(x - x_0)$ for some $x_0 \in \partial\Omega$, making $\text{data}_{\phi, \psi}$ the current at point x_0 when voltage ϕ is applied. With infinite pairs of ϕ and ψ in the experimental setup, EIT seeks to reconstruct $\bar{\sigma}(x)$. EIT further reduces to the famous Calderón problem when the span of ϕ and ψ covers the entire $H^{1/2}$. In practice, however, one typically has a rough estimate of the media $\bar{\sigma}(x)$, termed the background media $\sigma^*(x)$. (For example, most human lungs have the same structure.) In such situations, one can linearize and reconstruct the perturbation $\sigma(x) \stackrel{\text{def}}{=} \bar{\sigma}(x) - \sigma^*(x) \ll 1$. Specifically, suppose that ρ_1 solves the following background forward equation:

$$\begin{aligned}\nabla_x \cdot (\sigma^*(x) \nabla_x \rho_1(x)) &= 0, & x \in \Omega \\ \rho_1(x) &= \phi(x), & x \in \partial\Omega\end{aligned}\tag{14}$$

with the same boundary condition ϕ and the given known background media σ^* . Since both these quantities are known, $\rho_1(x)$ can be solved ahead of time for any ϕ . We can also define the adjoint equation:

$$\begin{cases} \nabla_x \cdot (\sigma^*(x) \nabla_x \rho_2(x)) = 0, & x \in \Omega \\ \rho_2(x) = \psi(x), & x \in \partial\Omega. \end{cases}\tag{15}$$

To obtain the Fredholm integral, we take the difference of (14) and (12) and drop higher order terms in $\sigma(x)$ to obtain

$$\begin{aligned}\nabla_x \cdot (\sigma^*(x) \nabla_x \rho(x)) &= -\nabla_x \cdot (\sigma(x) \nabla_x \rho_1(x)), & x \in \Omega \\ \rho(x) &= 0, & x \in \partial\Omega\end{aligned}\tag{16}$$

where $\rho(x) \stackrel{\text{def}}{=} \bar{\rho}(x) - \rho_1(x)$. With this equation multiplied with ρ_2 and the adjoint (15) multiplied with ρ , we integrate over Ω and integrating by parts. The left hand sides cancel and the right hand side of (16) will be balancing the boundary terms:

$$\int \nabla_x \rho_1(x) \cdot \nabla_x \rho_2(x) \sigma(x) dx = \int_{\partial\Omega} \sigma^* \frac{\partial \rho}{\partial n} \psi dx + \int_{\partial\Omega} \sigma \frac{\partial \rho_1}{\partial n} \psi dx.\tag{17}$$

While the left hand side of this equation is Fredholm integral testing on σ (the conductivity to be reconstructed) with test function $\nabla_x \rho_1(x) \cdot \nabla_x \rho_2(x)$, the right hand side is the data that we obtain from measurement $\text{data}_{\phi, \psi}$.

Indeed, since $\bar{\rho} = \rho + \rho_1$ and $\bar{\sigma} = \sigma + \sigma^*$, with $\rho \ll 1$ and $\sigma \ll 1$, the right hand side can be approximated by dropping the higher order term $\int_{\partial\Omega} \sigma \frac{\partial \rho}{\partial n} \psi dx$, as follows:

$$\begin{aligned} & \int_{\partial\Omega} \sigma^* \frac{\partial \rho}{\partial n} \psi dx + \int_{\partial\Omega} \sigma \frac{\partial \rho_1}{\partial n} \psi dx \\ &= \int_{\partial\Omega} \bar{\sigma}(x) \frac{\partial \bar{\rho}}{\partial n} \psi dx - \int_{\partial\Omega} \sigma^*(x) \frac{\partial \rho_1}{\partial n} \psi dx - \int_{\partial\Omega} \sigma \frac{\partial \rho}{\partial n} dx \\ &\approx \int_{\partial\Omega} \bar{\sigma}(x) \frac{\partial \bar{\rho}}{\partial n} \psi dx - \int_{\partial\Omega} \sigma^*(x) \frac{\partial \rho_1}{\partial n} \psi dx, \end{aligned}$$

This expression differs from $\overline{\text{data}}_{\phi, \psi}$ defined in (13) by $\int_{\partial\Omega} \sigma^* \frac{\partial \rho_1}{\partial n} \psi dx$, a pre-computed term, and thus the entire term is known. We finally have

$$\int \nabla_x \rho_1(x) \cdot \nabla_x \rho_2(x) \sigma(x) dx = \text{data}_{\phi, \psi}. \quad (18)$$

We emphasize that the ϕ dependence comes in through ρ_1 while the ψ dependence comes in through ρ_2 . These functions represent applied voltage source and measuring setup, respectively. If one can provide point source and point measurement, ϕ and ψ can be as sharp as Dirac-delta functions.

By varying ϕ and ψ , one finds infinitely many pairs $\{\rho_1(\cdot; \phi), \rho_2(\cdot; \psi)\}$, each pair providing one data point corresponding to one experiment setup. These experimental setup altogether give rise to an overdetermined Fredholm integral. More details can be found in [5, 7].

A similar problem arises in optical tomography [1]. Here we inject light into bio-tissue and take measurements of light intensity on the surface, to reconstruct the optical properties of the bio-tissue. The formulation is

$$\int \rho_1(x, v) \rho_2(x, v) \sigma(x, v) dx dv = \text{data}_{\phi, \psi}, \quad (19)$$

where $(x, v) \in \Omega \otimes \mathbb{S}$ (where Ω is the spatial domain and \mathbb{S} is the velocity domain), and ρ_i are solutions to the forward background radiative transfer equation and the adjoint equation:

$$\begin{cases} v \cdot \nabla_x \rho_1(x, v) = \sigma^*(x, v) \mathcal{L} \rho_1(x, v), & (x, v) \in \Omega \otimes \mathbb{S} \\ \rho_1(x, v) = 0, & (x, v) \in \Gamma_- \end{cases},$$

and

$$\begin{cases} -v \cdot \nabla_x \rho_2(x, v) = \sigma^*(x, v) \mathcal{L} \rho_2(x, v), & (x, v) \in \Omega \otimes \mathbb{S} \\ \rho_2(x, v) = \psi(x, v), & (x, v) \in \Gamma_+ \end{cases}.$$

In these equations, \mathcal{L} is a known integral linear operator on v , and Γ_- and Γ_+ are the set collecting incoming and outgoing boundary coordinates, namely $\Gamma_{\pm} = \{(x, v) : x \in \partial\Omega, \pm v \cdot n(x) > 0\}$ with $n(x)$ being an outer-normal direction at $x \in \partial\Omega$. By varying the boundary conditions ϕ and ψ , one can find infinitely many solution pairs of $\{\rho_1(\cdot, \phi), \rho_2(\cdot, \psi)\}$, and collect the corresponding data in (19). The inverse Fredholm integral (19) can then be solved for σ . We refer to [1, 6] for details of the linearization procedure.

When σ is discretized on p grid points, the reconstruction problem has the semi-infinite form $\mathbf{A}x \approx \mathbf{b}$, where $x \in \mathbb{R}^p$ is the discrete version of σ and \mathbf{A} and \mathbf{b} have infinitely many rows, corresponding to the infinitely many instances of ρ_1 and ρ_2 . A fully discrete version can be obtained by considering n_1 values of ρ_1 and n_2 values of ρ_2 , and setting $n = n_1 n_2$ to obtain a problem of the form (1). In the remainder of the paper, we study the sketched form of this system (2), for various choices of the sketching matrix \mathbf{S} .

3. SKETCHING WITH TENSOR STRUCTURES

We preface our results with a definition of (ε, δ) - l^2 embedding.

Definition 1 ((ε, δ) - l^2 embedding). Given matrix \bar{A} and $\varepsilon > 0$, let S be a random matrix drawn from a matrix distribution $(\Omega, \mathcal{F}, \Pi)$. If with probability at least $1 - \delta$, we have

$$|\|Sy\|^2 - \|y\|^2| \leq \varepsilon\|y\|^2, \quad \text{for all } y \in \text{Range}(\bar{A}), \quad (20)$$

then we say that S is an (ε, δ) - l^2 embedding of \bar{A} .

Note that (20) depends only on the space $\text{Range}(\bar{A})$ rather than the matrix itself, so we sometimes say instead that the random matrix S is an (ε, δ) - l^2 embedding of the linear vector space $\text{Range}(\bar{A})$. (We use the two terms interchangeably in discussions below.)

The (ε, δ) - l^2 embedding property is essentially the only property needed to bound the error resulting from sketching. It can be shown that if S is an (ε, δ) - l^2 embedding for the augmented matrix $\bar{A} \stackrel{\text{def}}{=} [A, b]$, then the two least-squares problems (1) and (2) are similar in the sense of (3), as the following result suggests.

Theorem 1. For $\varepsilon, \delta \in (0, 1/2)$, suppose that S is an (ε, δ) - l^2 embedding of the augmented matrix $\bar{A} \stackrel{\text{def}}{=} [A, b] \in \mathbb{R}^{n \times (p+1)}$. Then with probability at least $1 - \delta$, we have

$$\|Ax_s^* - b\|^2 \leq (1 + 4\varepsilon)\|Ax^* - b\|^2,$$

where x^* and x_s^* are defined in (1) and (2), respectively.

The proof of the theorem is rather standard. We simply use the definition of the (ε, δ) - l^2 embedding and the fact that:

$$(1 - \varepsilon)\|Ax_s^* - b\|^2 \leq \|S(Ax_s^* - b)\|^2 \leq \|S(Ax^* - b)\|^2 \leq (1 + \varepsilon)\|Ax^* - b\|^2.$$

For $0 \leq \varepsilon \leq 1/2$, this leads to

$$\|Ax_s^* - b\|^2 \leq \frac{1 + \varepsilon}{1 - \varepsilon}\|Ax^* - b\|^2 \leq (1 + 4\varepsilon)\|Ax^* - b\|^2.$$

Given this result, we focus henceforth on whether the various sampling strategies form an (ε, δ) - l^2 embedding of the augmented matrix $\bar{A} = [A, b]$.

Another theorem that is crucial to our analysis, proved in [40], states that Gaussian matrices are (ε, δ) - l^2 embeddings if the number of rows is sufficiently large. This result does not consider tensor structure of A .

Theorem 2 (Theorem 2.3 from [40]). Let $R \in \mathbb{R}^{r \times n}$ be a Gaussian matrix, meaning that each entry R_{ij} is drawn i.i.d. from a normal distribution $\mathcal{N}(0, 1)$, and define $S \in \mathbb{R}^{r \times n}$ to be the scaled Gaussian matrix defined by

$$S = \frac{1}{\sqrt{r}}R.$$

For any fixed matrix $A \in \mathbb{R}^{n \times p}$ and $\varepsilon, \delta \in (0, 1/2)$, this choice of S is an (ε, δ) - l^2 embedding of A provided that

$$r \geq \frac{C}{\varepsilon^2}(|\log \delta| + p),$$

where $C > 0$ is a constant independent of ε, δ, n , and p .

The lower bound of r is almost optimal for the sketched regression problem: the bound is independent of the number of equations n , and grows only linearly in the number of unknowns p . That is, the numbers of equations and unknowns in the sketched problem (2) are of the same order. The theorem is proved by constructing a γ -net for the unit sphere in $\text{Range}(A)$ and applying the Johnson-Lindenstrauss lemma.

Building on the concept of (ε, δ) - l^2 embedding and the relationship between (ε, δ) - l^2 embedding and sketching (Theorem 1), we will study the lower bound for r (the number of rows needed in the sketching) when the tensor

structure of Case 1 or Case 2 is imposed. Our basic strategy is to decompose the tensor structure into smaller components to which Theorem 2 can be applied.

We state the results below and present proofs in Sections 4 and 5 for the two different cases.

Recall the notation that we defined in Section 1. The matrices F , G are defined in (5) and A is defined in (6). Both F and G are assumed to have full column rank p . We need to design the sketching matrix S to (ε, δ) - l^2 embed $\text{Range}(\bar{A})$, the space spanned by $\{f_b \otimes g_b\} \cup \{a_j \stackrel{\text{def}}{=} f_j \otimes g_j, j = 1, \dots, p\}$. In Theorem 3 and 4, we construct the (ε, δ) - l^2 embedding matrix of the Kronecker product $F \otimes G$, which automatically becomes a (ε, δ) - l^2 embedding of its column submatrix A . Moreover, we show in Corollaries 1 and 2 that these results can be extended to construct (ε, δ) - l^2 embeddings of the augmented matrix \bar{A} by constructing (ε, δ) - l^2 embeddings of the Kronecker product of the augmented matrices $\bar{F} \otimes \bar{G}$, where

$$\bar{F} = [F, f_b], \quad \bar{G} = [G, g_b]. \quad (21)$$

For Case 1, we have the following result.

Theorem 3. Consider $S = P \otimes Q \in \mathbb{R}^{r_1 r_2 \times n_1 n_2}$ where $P \in \mathbb{R}^{r_1 \times n_1}$, $Q \in \mathbb{R}^{r_2 \times n_2}$ are independent scaled Gaussian matrices defined by

$$P \stackrel{\text{def}}{=} \frac{1}{\sqrt{r_1}} R \quad \text{and} \quad Q \stackrel{\text{def}}{=} \frac{1}{\sqrt{r_2}} R', \quad \text{where } R_{ij}, R'_{ij} \text{ are i.i.d. normal for all } i, j.$$

For any given full rank matrices $F \in \mathbb{R}^{n_1 \times p}$, $G \in \mathbb{R}^{n_2 \times p}$, and $A \in \mathbb{R}^{n \times p}$ as in (5) and (6), and $\varepsilon, \delta \in (0, 1/2)$, the random matrix S is an (ε, δ) - l^2 embedding of $F \otimes G$ and A provided that

$$r_i \geq \frac{C}{\varepsilon^2} (|\log \delta| + p), \quad i = 1, 2, \quad (22)$$

where the constant $C > 0$ is independent of $\varepsilon, \delta, n_1, n_2$, and p .

Corollary 1. Consider the matrices S , F , G , and A from Theorem 3, and assume that the vector b has the form (7). Then for given $\varepsilon, \delta \in (0, 1/2)$, the random matrix S is an (ε, δ) - l^2 embedding of the augmented matrix $\bar{A} \stackrel{\text{def}}{=} [A, b]$, provided that

$$r_i \geq \frac{C}{\varepsilon^2} (|\log \delta| + p + 1), \quad i = 1, 2, \quad (23)$$

where the constant $C > 0$ is independent of $\varepsilon, \delta, n_1, n_2$, and p .

Proof. Define the augmented matrices \bar{F} and \bar{G} as in (21). We have that

$$\text{Range}(\bar{F} \otimes \bar{G}) = \text{Span}\{F \otimes G, f_1 \otimes g_b, \dots, f_p \otimes g_b, f_b \otimes g_1, \dots, f_b \otimes g_p, b\}.$$

Supposing that \bar{F} and \bar{G} have full rank, the linear subspace $\text{Range}(\bar{A})$ is a subspace of $\text{Range}(\bar{F} \otimes \bar{G})$. By applying Theorem 3 to the augmented matrices \bar{F} and \bar{G} and using (23), we have that S is an (ε, δ) - l^2 embedding of $\text{Range}(\bar{F} \otimes \bar{G})$ as well as its subspace $\text{Range}(\bar{A})$. Supposing that \bar{F} is not of full rank but \bar{G} is of full rank, the subspace $\text{Range}(\bar{F} \otimes \bar{G})$ is a subspace of $\text{Range}(F \otimes \bar{G})$, so similar results can be obtained by applying Theorem 3 to F and \bar{G} . Other cases regarding the rank of \bar{F} and \bar{G} can be dealt with in the same way. \square

The result for Case 2 is as follows.

Theorem 4. Let $p_i \in \mathbb{R}^{n_1}$, $q_i \in \mathbb{R}^{n_2}$, $i = 1, 2, \dots, r$ be independent random Gaussian vectors, and define the sketching matrix S to have the form:

$$S = \frac{1}{\sqrt{r}} \begin{bmatrix} p_1^\top \otimes q_1^\top \\ \vdots \\ p_r^\top \otimes q_r^\top \end{bmatrix} \in \mathbb{R}^{r \times n_1 n_2}. \quad (24)$$

Suppose that $p \geq 6$, and that $\mathbf{F} \in \mathbb{R}^{n_1 \times p}$, $\mathbf{G} \in \mathbb{R}^{n_2 \times p}$, and $\mathbf{A} \in \mathbb{R}^{n \times p}$ are full-rank matrices defined as in (5) and (6). Let $\varepsilon, \delta \in (0, 1/2)$. Then the random matrix \mathbf{S} is an (ε, δ) - l^2 embedding of $\mathbf{F} \otimes \mathbf{G}$ and \mathbf{A} provided that

$$r \geq C \max \left\{ \frac{1}{\varepsilon} (|\log \delta| + p^2)^3, \frac{1}{\varepsilon^{5/2}} \right\}, \quad (25)$$

where $C > 0$ is a constant independent of ε , δ , n_1 , n_2 , and p .

Corollary 2. Consider the same matrices \mathbf{S} , \mathbf{F} , \mathbf{G} , and \mathbf{A} as in Theorem 4, with $p \geq 6$, and assume that vector \mathbf{b} is of the form (7). Then for given $\varepsilon, \delta \in (0, 1/2)$, the random matrix \mathbf{S} is an (ε, δ) - l^2 embedding of the augmented matrix $\bar{\mathbf{A}} \stackrel{\text{def}}{=} [\mathbf{A}, \mathbf{b}]$ provided that

$$r \geq C \max \left\{ \frac{1}{\varepsilon} (|\log \delta| + (p+1)^2)^3, \frac{1}{\varepsilon^{5/2}} \right\}, \quad (26)$$

where the constant $C > 0$ is independent of ε , δ , n_1 , n_2 , and p .

We omit the proof since it is similar to that of Corollary 1.

Theorems 1 and 2 yield the fundamental results that, with high probability, for any fixed overdetermined linear problem, the sketched problem in which \mathbf{S} is a Gaussian matrix can achieve optimal residual up to a small multiplicative error. In particular, as will be clear in the proof later, the Case-1 tensor-structured sketching matrix $\mathbf{S} = \mathbf{P} \otimes \mathbf{Q}$ not only (ε, δ) - l^2 embeds $\mathbf{A} = \mathbf{F} \otimes \mathbf{G}$, but the number of rows in \mathbf{P} and \mathbf{Q} each depends only linearly on p (see (22)), so that the number of rows in \mathbf{S} scales like p^2 . If the Case-2 sketching matrix is used, the dependence of r on p and ε is more complex. Whether this bound is greater than or less than the bound for Case 1 depends on the relative sizes of ε^{-1} and p .

We stress that both bounds show that the number of rows in \mathbf{S} is independent of the dimension $n \stackrel{\text{def}}{=} n_1 n_2$ of the ambient space. This allows n to be potentially infinity. We also stress that the dependence on ε and p may not be optimal, and the bound may not be tight. As will be seen in the later sections, we have limited understanding of quartic powers of Gaussian random variables, and this confines us obtaining a tighter bound.

4. CASE 1: PROOF OF THEOREM 3

In this section we present the proof of Theorem 3. We start with technical results.

Lemma 1. Consider natural numbers r_2 , n_1 , and n_2 , and assume that a random matrix $\mathbf{Q} \in \mathbb{R}^{r_2 \times n_2}$ is an (ε, δ) - l^2 embedding of \mathbb{R}^{n_2} , meaning that with probability at least $1 - \delta$, \mathbf{Q} preserves l^2 norm with ε accuracy, that is,

$$|\|\mathbf{Q}\mathbf{x}\|^2 - \|\mathbf{x}\|^2| \leq \varepsilon \|\mathbf{x}\|^2, \quad \text{for all } \mathbf{x} \in \mathbb{R}^{n_2}.$$

Then the Kronecker product $\text{Id}_{n_1} \otimes \mathbf{Q}$ is an (ε, δ) - l^2 embedding of $\mathbb{R}^{n_1 n_2}$. Similarly, if $\mathbf{Q} \in \mathbb{R}^{r_1 \times n_1}$ is an (ε, δ) - l^2 embedding of \mathbb{R}^{n_1} , then $\mathbf{Q} \otimes \text{Id}_{n_2}$ is an (ε, δ) - l^2 embedding of $\mathbb{R}^{n_1 n_2}$.

Proof. The proof for the two statements are rather similar, so we prove only the first claim.

Any $\mathbf{x} \in \mathbb{R}^{n_1 n_2}$ can be written in the following form

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{n_1} \end{bmatrix}, \quad \text{where } \mathbf{x}_i \in \mathbb{R}^{n_2}, i = 1, 2, \dots, n_1.$$

Then

$$(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x} = \begin{bmatrix} \mathbf{Q} & & \\ & \ddots & \\ & & \mathbf{Q} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{n_1} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}\mathbf{x}_1 \\ \vdots \\ \mathbf{Q}\mathbf{x}_{n_1} \end{bmatrix}.$$

Thus, we have

$$\|(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}\|^2 = \sum_{i=1}^n \|\mathbf{Q}\mathbf{x}_i\|^2, \quad \|\mathbf{x}\|^2 = \sum_{i=1}^n \|\mathbf{x}_i\|^2. \quad (27)$$

Since \mathbf{Q} is an (ε, δ) - l^2 embedding of \mathbb{R}^{n_2} , then with probability at least $1 - \delta$, for all $\mathbf{x}_i \in \mathbb{R}^{n_2}$, we have

$$|\|\mathbf{Q}\mathbf{x}_i\|^2 - \|\mathbf{x}_i\|^2| \leq \varepsilon \|\mathbf{x}_i\|^2, \quad \text{for all } i = 1, 2, \dots, n_1. \quad (28)$$

By using this bound in (27), with probability at least $1 - \delta$, we have for all $\mathbf{x} \in \mathbb{R}^{n_1 n_2}$ that

$$|\|(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}\|^2 - \|\mathbf{x}\|^2| \leq \sum_{i=1}^n |\|\mathbf{Q}\mathbf{x}_i\|^2 - \|\mathbf{x}_i\|^2| \leq \varepsilon \sum_{i=1}^n \|\mathbf{x}_i\|^2 = \varepsilon \|\mathbf{x}\|^2,$$

so that $(\text{Id}_{n_1} \otimes \mathbf{Q})$ is an (ε, δ) - l^2 embedding of $\mathbb{R}^{n_1 n_2}$, as claimed. \square

The following corollary extends the previous result and discusses the embedding property of $\mathbf{P} \otimes \mathbf{Q}$.

Corollary 3. Assume two random matrices $\mathbf{P} \in \mathbb{R}^{r_1 \times n_1}$ and $\mathbf{Q} \in \mathbb{R}^{r_2 \times n_2}$ are (ε, δ) - l^2 embeddings of \mathbb{R}^{n_1} and \mathbb{R}^{n_2} , respectively. Then the Kronecker product $\mathbf{P} \otimes \mathbf{Q} \in \mathbb{R}^{r_1 r_2 \times n_1 n_2}$ is an $(\varepsilon(2 + \varepsilon), 2\delta)$ - l^2 embedding of $\mathbb{R}^{n_1 n_2}$.

Proof. Noting that (see (68) in Appendix A),

$$\mathbf{P} \otimes \mathbf{Q} = (\mathbf{P} \otimes \text{Id}_{r_2})(\text{Id}_{n_1} \otimes \mathbf{Q}),$$

we have

$$\|(\mathbf{P} \otimes \mathbf{Q})\mathbf{x}\|^2 = \|(\mathbf{P} \otimes \text{Id}_{r_2})(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}\|^2 = \|(\mathbf{P} \otimes \text{Id}_{r_2})\mathbf{y}\|^2,$$

where $\mathbf{y} \stackrel{\text{def}}{=} (\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}$.

Denote by $(\Omega_1, \mathcal{F}_1, \Pi_1)$ and $(\Omega_2, \mathcal{F}_2, \Pi_2)$ the probability triplets for \mathbf{P} and \mathbf{Q} , respectively. Since \mathbf{P} is an (ε, δ) - l^2 embedding of \mathbb{R}^{n_1} , we have with probability at least $1 - \delta$ in Π_1 that

$$|\|(\mathbf{P} \otimes \text{Id}_{r_2})\mathbf{y}\|^2 - \|\mathbf{y}\|^2| \leq \varepsilon \|\mathbf{y}\|^2.$$

Similarly, with probability at least $1 - \delta$ for the choice of \mathbf{Q} in Π_2 , we have

$$|\|(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}\|^2 - \|\mathbf{x}\|^2| \leq \varepsilon \|\mathbf{x}\|^2, \quad \text{for all } \mathbf{x} \in \mathbb{R}^{n_1 n_2}.$$

Combining the two inequalities, we have with probability at least $1 - 2\delta$ in the joint probability space of Π_1 and Π_2 that the following is true for all $\mathbf{x} \in \mathbb{R}^{n_1 n_2}$:

$$\begin{aligned} |\|(\mathbf{P} \otimes \mathbf{Q})\mathbf{x}\|^2 - \|\mathbf{x}\|^2| &\leq |\|(\mathbf{P} \otimes \text{Id}_{r_2})\mathbf{y}\|^2 - \|\mathbf{y}\|^2| + |\|(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}\|^2 - \|\mathbf{x}\|^2| \\ &\leq \varepsilon \|\mathbf{y}\|^2 + \varepsilon \|\mathbf{x}\|^2 \\ &= \varepsilon \|(\text{Id}_{n_1} \otimes \mathbf{Q})\mathbf{x}\|^2 + \varepsilon \|\mathbf{x}\|^2 \\ &\leq \varepsilon(2 + \varepsilon) \|\mathbf{x}\|^2. \end{aligned}$$

This concludes the proof. \square

Now we are ready to show the proof of Theorem 3, obtained by applying Theorem 2 to Corollary 3.

Proof of Theorem 3. For any vector \mathbf{y} in the span of $\mathbf{F} \otimes \mathbf{G}$, we can write

$$\mathbf{y} = (\mathbf{U}_F \otimes \mathbf{U}_G)\mathbf{x}, \quad \text{for some } \mathbf{x} \in \mathbb{R}^{p^2},$$

where $\mathbf{U}_F \in \mathbb{R}^{n_1 \times p}$ and $\mathbf{U}_G \in \mathbb{R}^{n_2 \times p}$ collect the left singular vectors of matrices \mathbf{F} and \mathbf{G} , respectively. By applying (68) from Appendix A, we have

$$(\mathbf{U}_F \otimes \mathbf{U}_G) = (\mathbf{U}_F \otimes \text{Id}_{n_2})(\text{Id}_p \otimes \mathbf{U}_G).$$

It is easy to see that the matrix $\text{Id}_p \otimes \mathbf{U}_G$ has orthonormal columns, so it is an isometry. The matrices $\mathbf{U}_F \otimes \text{Id}_{n_2}$ and $\mathbf{U}_F \otimes \mathbf{U}_G$ are isometries for the same reason. As a consequence, we have $\|y\|^2 = \|x\|^2$. From (68) in Appendix A, we have by defining $\tilde{\mathbf{P}} \stackrel{\text{def}}{=} \mathbf{P}\mathbf{U}_F \in \mathbb{R}^{r_1 \times p}$ and $\tilde{\mathbf{Q}} \stackrel{\text{def}}{=} \mathbf{Q}\mathbf{U}_G \in \mathbb{R}^{r_2 \times p}$ that

$$\mathbf{S}y = (\mathbf{P} \otimes \mathbf{Q})(\mathbf{U}_F \otimes \mathbf{U}_G)x = (\mathbf{P}\mathbf{U}_F) \otimes (\mathbf{Q}\mathbf{U}_G)x = (\tilde{\mathbf{P}} \otimes \tilde{\mathbf{Q}})x. \quad (29)$$

Due to the orthogonality of \mathbf{U}_F and \mathbf{U}_G , the random matrices $\tilde{\mathbf{P}}$ and $\tilde{\mathbf{Q}}$ are also independent Gaussian matrices with i.i.d. entries. According to Theorem 2, for any pair $\tilde{\varepsilon}, \tilde{\delta} \in (0, 1/2)$, by choosing r_i to satisfy

$$r_i \geq \frac{C}{\tilde{\varepsilon}^2}(|\log \tilde{\delta}| + p), \quad i = 1, 2, \quad (30)$$

we have that $\tilde{\mathbf{P}}$ and $\tilde{\mathbf{Q}}$ are both $(\tilde{\varepsilon}, \tilde{\delta})$ - l^2 embeddings of \mathbb{R}^p . Thus, from Corollary 3, the tensor product $(\tilde{\mathbf{P}} \otimes \tilde{\mathbf{Q}})$ is an $(\tilde{\varepsilon}(2 + \tilde{\varepsilon}), 2\tilde{\delta})$ - l^2 embedding of \mathbb{R}^{p^2} , meaning that with probability at least $1 - 2\tilde{\delta}$, we have

$$\left| \|(\tilde{\mathbf{P}} \otimes \tilde{\mathbf{Q}})x\|^2 - \|x\|^2 \right| \leq \tilde{\varepsilon}(2 + \tilde{\varepsilon})\|x\|^2, \quad \text{for all } x \in \mathbb{R}^{p^2}.$$

Recalling $\|x\|^2 = \|y\|^2$ and (29), we have that

$$\left| \|\mathbf{S}y\|^2 - \|y\|^2 \right| \leq \tilde{\varepsilon}(2 + \tilde{\varepsilon})\|y\|^2, \quad \text{for all } y \in \text{Span}\{\mathbf{F} \otimes \mathbf{G}\}.$$

By defining $\varepsilon = \tilde{\varepsilon}(2 + \tilde{\varepsilon})$ and $\delta = 2\tilde{\delta}$, we have

$$\tilde{\varepsilon} = \frac{\varepsilon}{\sqrt{1 + \varepsilon} + 1}, \quad \text{and} \quad \tilde{\delta} = \frac{\delta}{2}.$$

Note that if ε and δ are in $(0, 1/2)$, then $\tilde{\varepsilon}$ and $\tilde{\delta}$ are also in this interval, so (30) applies. By substituting into (30) we obtain

$$r_i \geq \frac{C}{\varepsilon^2}(|\log \delta| + p), \quad i = 1, 2.$$

The constant C here is different from the value in (30) but can still be chosen independently of $\varepsilon, \delta, n_1, n_2$, and p . We conclude that $\mathbf{S} = \mathbf{P} \otimes \mathbf{Q}$ is an (ε, δ) - l^2 embedding of $\mathbf{F} \otimes \mathbf{G}$ and thus also an (ε, δ) - l^2 embedding of \mathbf{A} . \square

5. CASE 2: PROOF OF THEOREM 4

In this section we investigate Case-2 sketching matrices, which have the form (24).

We prove Theorem 4 in two major steps. First, in Section 5.1, we investigate the accuracy and probability of embedding any given vector $y \in \text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$. Second, in Section 5.2, we extend this study to deal with the whole space $\text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$. To do so, we first build a γ -net over the unit sphere in $\text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$ so that we can “approximate” the space using a finite set of vectors. By adjusting ε and δ , one not only preserves the norm, but also the angles between the vectors on the net. We then map the net back to the space to show that \mathbf{S} preserves the norm of the vectors in the whole space. This standard technique is used in [40] to prove their Theorem 2.

5.1. Embedding a given vector. We establish the following result, whose proof appears at the end of the subsection.

Proposition 1. *Given two full rank matrices \mathbf{F} and \mathbf{G} as in (5) and $\varepsilon \in (0, 1/2)$, let $\mathbf{S} \in \mathbb{R}^{r \times n_1 n_2}$ have the form of (24), with \mathbf{p}_i and \mathbf{q}_i , $i = 1, 2, \dots, r$ being i.i.d. Gaussian vectors. Then for any fixed $y \in \text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$, we have that*

$$\Pr(|\|\mathbf{S}y\|^2 - \|y\|^2| > \varepsilon\|y\|^2) \leq 5r \exp\left(\frac{3}{4}p^{1/2}\right) \exp\left(-\frac{1}{2}r^{1/3}\varepsilon^{1/3}\right),$$

provided that

$$r \geq 8 \cdot 3^{3/2} \cdot \max\{\varepsilon^{-5/2}, p^{3/2}\varepsilon^{-1}\}.$$

Essentially, this proposition says that \mathbf{S} is an $(\varepsilon, 5r \exp((3/4)p^{1/2}) \exp(-(1/2)r^{1/3}\varepsilon^{1/3}))$ - l^2 embedding of any fixed $\mathbf{y} \in \text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$. The contribution from the factor $\exp(-(1/2)r^{1/3}\varepsilon^{1/3})$ is small when r is large.

We start with several technical lemmas. Lemma 2 identifies $|\|\mathbf{S}\mathbf{y}\|^2 - \|\mathbf{y}\|^2|/\|\mathbf{y}\|$ with a particular type of random variable; we discuss the tail bound for this random variable in Lemma 4. Lemma 3 contains some crucial estimates to be used in Lemma 4.

Lemma 2. *Given two full rank matrices \mathbf{F} and \mathbf{G} as in (5), consider \mathbf{S} defined as in (24). Then there exists a diagonal positive semi-definite matrix Σ with $\text{Tr}(\Sigma^2) = 1$ so that for any $\mathbf{y} \in \text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$ with $\|\mathbf{y}\| = 1$, we have*

$$\|\mathbf{S}\mathbf{y}\|^2 \stackrel{d}{\sim} \frac{1}{r} \sum_{i=1}^r \zeta_i^2, \quad \text{where } \zeta_i \stackrel{\text{def}}{=} \xi_i^\top \Sigma \eta_i,$$

where $\stackrel{d}{\sim}$ denotes equal in distribution and $\xi_i, \eta_i \in \mathbb{R}^p$ are independent Gaussian vectors drawn from $\mathcal{N}(0, \text{Id}_p)$.

Proof. From (24) we have

$$\mathbf{S}\mathbf{y} = \frac{1}{\sqrt{r}} \begin{bmatrix} (\mathbf{p}_1^\top \otimes \mathbf{q}_1^\top) \mathbf{y} \\ \vdots \\ (\mathbf{p}_r^\top \otimes \mathbf{q}_r^\top) \mathbf{y} \end{bmatrix} \implies \|\mathbf{S}\mathbf{y}\|^2 = \frac{1}{r} \sum_{i=1}^r \zeta_i^2,$$

where $\zeta_i \stackrel{\text{def}}{=} (\mathbf{p}_i^\top \otimes \mathbf{q}_i^\top) \mathbf{y}$. Since \mathbf{p}_i and \mathbf{q}_i are independent Gaussian vectors, all random variables ζ_i , $i = 1, 2, \dots, r$, are drawn i.i.d. from the same distribution.

We consider now the behavior of $\zeta \stackrel{\text{def}}{=} (\mathbf{p}^\top \otimes \mathbf{q}^\top) \mathbf{y}$ for Gaussian vectors \mathbf{p} and \mathbf{q} . Notice that for any $\mathbf{y} \in \mathbb{R}^{n_1 n_2} \in \text{Span}\{\mathbf{F} \otimes \mathbf{G}\}$, there exists $\mathbf{x} \in \mathbb{R}^{p^2}$ such that

$$\mathbf{y} = (\mathbf{U}_\mathbf{F} \otimes \mathbf{U}_\mathbf{G}) \mathbf{x}, \quad \text{with } \|\mathbf{x}\| = 1,$$

where $\mathbf{U}_\mathbf{F}$ and $\mathbf{U}_\mathbf{G}$ collect the left singular vectors of \mathbf{F} and \mathbf{G} , respectively. We thus obtain from (68) that

$$\zeta = (\mathbf{p}^\top \otimes \mathbf{q}^\top) \mathbf{y} = (\mathbf{p}^\top \otimes \mathbf{q}^\top) (\mathbf{U}_\mathbf{F} \otimes \mathbf{U}_\mathbf{G}) \mathbf{x} = ((\mathbf{p}^\top \mathbf{U}_\mathbf{F}) \otimes (\mathbf{q}^\top \mathbf{U}_\mathbf{G})) \mathbf{x} = (\tilde{\mathbf{p}}^\top \otimes \tilde{\mathbf{q}}^\top) \mathbf{x},$$

where $\tilde{\mathbf{p}} \stackrel{\text{def}}{=} \mathbf{U}_\mathbf{F}^\top \mathbf{p} \in \mathbb{R}^p$ and $\tilde{\mathbf{q}} \stackrel{\text{def}}{=} \mathbf{U}_\mathbf{G}^\top \mathbf{q} \in \mathbb{R}^p$ are i.i.d. Gaussian vectors as well. By applying (69) and (70), we obtain

$$\zeta = (\tilde{\mathbf{p}}^\top \otimes \tilde{\mathbf{q}}^\top) \mathbf{x} = \tilde{\mathbf{q}}^\top \mathbf{Mat}(\mathbf{x}) \tilde{\mathbf{p}}, \quad (31)$$

where $\mathbf{Mat}(\mathbf{x}) \in \mathbb{R}^{p \times p}$ is the matricization of \mathbf{x} , discussed in Appendix A. By using the singular value decomposition $\mathbf{Mat}(\mathbf{x}) = \mathbf{U} \Sigma \mathbf{V}^\top$, we obtain

$$\text{Tr}(\Sigma^2) = \|\mathbf{Mat}(\mathbf{x})\|_F^2 = \|\mathbf{x}\|^2 = 1,$$

where $\|\cdot\|_F$ denotes the Frobenius norm of a matrix. By substituting into (31), we obtain

$$\zeta = (\mathbf{U}^\top \tilde{\mathbf{q}})^\top \Sigma \mathbf{V}^\top \tilde{\mathbf{p}} = \xi^\top \Sigma \eta,$$

where

$$\xi \stackrel{\text{def}}{=} \mathbf{U}^\top \tilde{\mathbf{q}} \in \mathbb{R}^p \quad \text{and} \quad \eta \stackrel{\text{def}}{=} \mathbf{V}^\top \tilde{\mathbf{p}} \in \mathbb{R}^p$$

are again i.i.d. Gaussian vectors in \mathbb{R}^p . This completes the proof. \square

Lemma 3. *For any fixed diagonal semi-positive definite matrix $\Sigma \stackrel{\text{def}}{=} \text{diag}\{\sigma_1, \dots, \sigma_p\}$ such that $\text{Tr}(\Sigma^2) = 1$, define the random variable ζ to be $\zeta \stackrel{\text{def}}{=} \xi^\top \Sigma \eta$, with ξ and η being i.i.d. random Gaussian vectors with p components. Then ζ satisfies the following properties:*

1.

$$\Pr(|\zeta| > t) \leq \begin{cases} 2 \exp\left(-\frac{(t-\sqrt{p})^2}{4\sqrt{p}}\right) & \text{if } \sqrt{p} \leq t \leq 2\sqrt{p} \\ 2 \exp\left(-\frac{(2t-3\sqrt{p})}{4}\right) & \text{if } t \geq 2\sqrt{p} \end{cases}, \quad (32)$$

2.

$$\mathbb{E}[\zeta^2] = 1 \quad \text{and} \quad \mathbb{E}[\zeta^4] \leq 9, \quad (33)$$

3.

$$\mathbb{E}[(|\zeta|^2 - \mathbb{E}[\zeta^2])^2] \leq 8. \quad (34)$$

Proof. For any $s > 0$ and $t \geq 0$, we apply Markov's inequality to derive

$$\Pr(\zeta > t) = \Pr(e^{s\zeta} > e^{st}) \leq e^{-st} \mathbb{E}[\exp(s\xi^\top \Sigma \eta)]. \quad (35)$$

Noting that $2\xi^\top \Sigma \eta \leq \|\Sigma^{1/2}\xi\|^2 + \|\Sigma^{1/2}\eta\|^2$, we use the independence of ξ and η to deduce that

$$\mathbb{E}[\exp(s\xi^\top \Sigma \eta)] \leq \mathbb{E}\left[\exp\left((s/2)(\|\Sigma^{1/2}\xi\|^2 + \|\Sigma^{1/2}\eta\|^2)\right)\right] = \mathbb{E}\left[e^{(s/2)\|\Sigma^{1/2}\xi\|^2}\right] \mathbb{E}\left[e^{(s/2)\|\Sigma^{1/2}\eta\|^2}\right]. \quad (36)$$

For the first term on the right-hand side of (36), using independence of the ξ_i and the concave Jensen's inequality, we have that

$$\mathbb{E}\left[e^{(s/2)\|\Sigma^{1/2}\xi\|^2}\right] = \mathbb{E}\left[\exp\left(\frac{s}{2} \sum_{i=1}^p \sigma_i \xi_i^2\right)\right] = \prod_{i=1}^p \mathbb{E}\left[e^{s\xi_i^2/2}\right] \leq \prod_{i=1}^p \left(\mathbb{E}\left[e^{s\xi_i^2/2}\right]\right)^{\sigma_i},$$

where we used $0 \leq \sigma_i \leq 1$, $i = 1, 2, \dots, p$ to apply the concave Jensen's inequality, and $\xi_i \sim \mathcal{N}(0, 1)$. According to Proposition 2 (see Appendix A.2), $\xi_i^2 - 1$ is a sub-exponential random variable with parameters $(2, 4)$. Thus from (71), with $\lambda = 2$, $b = 4$, and s replaced by $s/2$, we have

$$\begin{aligned} \mathbb{E}\left[e^{(s/2)\|\Sigma^{1/2}\xi\|^2}\right] &\leq \prod_{i=1}^p \left(\mathbb{E}_\xi\left[e^{s\xi_i^2/2}\right]\right)^{\sigma_i} = \left(\mathbb{E}_\xi\left[e^{s\xi^2/2}\right]\right)^{\text{Tr}(\Sigma)} \\ &= \left(e^{s/2} \mathbb{E}_\xi\left[e^{s(\xi^2-1)/2}\right]\right)^{\text{Tr}(\Sigma)} \leq e^{(s^2+s)\text{Tr}(\Sigma)/2}, \quad \text{for } 0 < s < 1/2. \end{aligned}$$

Since, by Hölder's inequality, we have

$$\text{Tr}(\Sigma) = \sum_{i=1}^p \sigma_i \leq \left(\sum_{i=1}^p \sigma_i^2\right)^{1/2} \sqrt{p} = \sqrt{p},$$

it follows that

$$\mathbb{E}\left[e^{(s/2)\|\Sigma^{1/2}\xi\|^2}\right] \leq e^{(s^2+s)\sqrt{p}/2}, \quad \text{for } s \in (0, 1/2).$$

The same bound holds for second term on the right-hand side of (36). When we substitute these bounds into (35) and (36), we obtain

$$\Pr(\zeta > t) \leq \exp(\sqrt{p}s^2 - (t - \sqrt{p})s).$$

By minimizing the right-hand side over $s \in [0, 1/2]$, we obtain

$$\Pr(\zeta > t) \leq \begin{cases} e^{-\frac{(t-\sqrt{p})^2}{4\sqrt{p}}} & \text{if } \sqrt{p} \leq t \leq 2\sqrt{p} \\ e^{-\frac{(2t-3\sqrt{p})}{4}} & \text{if } t \geq 2\sqrt{p} \end{cases}.$$

Due to symmetry, we have the same bound for $\Pr(\zeta < -t)$, so (32) follows.

To show the second statement, we notice that

$$\mathbb{E}[(\zeta^2 - \mathbb{E}[\zeta^2])^2] = \mathbb{E}[\zeta^4] - (\mathbb{E}[\zeta^2])^2. \quad (37)$$

By considering $\zeta = \sum_{i=1}^p \sigma_i \xi_i \eta_i$, the second moment can be calculated directly:

$$\mathbb{E}[\zeta^2] = \mathbb{E}\left[\sum_{i,j=1}^p \sigma_i \sigma_j \xi_i \xi_j \eta_i \eta_j\right] = \mathbb{E}\left[\sum_{i=1}^p \sigma_i^2 \xi_i^2 \eta_i^2\right] = \sum_{i=1}^p \sigma_i^2 = 1, \quad (38)$$

where we used the independence of ξ_i and η_i , the fact that $\mathbb{E}\xi_i = \mathbb{E}\eta_i = 0$ and $\mathbb{E}\xi_i^2 = \mathbb{E}\eta_i^2 = 1$.

To control the fourth moment, we notice that

$$\mathbb{E}[\zeta^4] = \mathbb{E}\left[\sum_{i,j,k,l} \sigma_i \sigma_j \sigma_k \sigma_l \xi_i \xi_j \xi_k \xi_l \eta_i \eta_j \eta_k \eta_l\right].$$

Due to the independence and the fact that all odd moments vanish for Gaussian random variables, the only terms in the summation that survive either have all indices equal ($i = j = l = k$) or two indices equal to one value while the other two indices equal a different value, for example $i = j$ and $k = l$ but $i \neq k$. Altogether, we obtain

$$\mathbb{E}[\zeta^4] = 3\mathbb{E}\left[\sum_{i \neq k} \sigma_i^2 \sigma_k^2 \xi_i^2 \xi_k^2 \eta_i^2 \eta_k^2\right] + \mathbb{E}\left[\sum_i \sigma_i^4 \xi_i^4 \eta_i^4\right],$$

where the coefficient in front of the first term comes from $\binom{4}{2}/\binom{2}{1} = 3$. Considering $\mathbb{E}\xi^2 = 1$ and $\mathbb{E}\xi^4 = 3$, we have

$$\begin{aligned} \mathbb{E}[\zeta^4] &= 3 \sum_{i \neq k} \sigma_i^2 \sigma_k^2 + 9 \sum_i \sigma_i^4 = 3 \sum_{i,k=1}^p \sigma_i^2 \sigma_k^2 + 6 \sum_i \sigma_i^4 \\ &\leq 3 \sum_{i,k=1}^p \sigma_i^2 \sigma_k^2 + 6 \sum_i \sigma_i^2 = 3 \left(\sum_{i=1}^p \sigma_i^2\right) \left(\sum_{k=1}^p \sigma_k^2\right) + 6 \sum_{i=1}^p \sigma_i^2 = 9, \end{aligned} \quad (39)$$

where we used $\sigma_i^4 \leq \sigma_i^2$. By substituting (38) and (39) into (37), we have

$$\mathbb{E}\left[\left(|\zeta|^2 - \mathbb{E}[\zeta^2]\right)^2\right] = \mathbb{E}[\zeta^4] - (\mathbb{E}[\zeta^2])^2 \leq 9 - 1^2 = 8,$$

which concludes the proof. \square

Remark 1. We note that this lemma is not new; its proof can be made more compact if one uses Hanson-Wright inequality and [37, Lemma 6.2.2]. The latter result shows that there exist absolute positive constants c and C such that

$$\mathbb{E}[\exp(s\xi^\top \Sigma \eta)] \leq \exp(Cs^2)$$

for $|s| \leq \frac{c}{\sigma_1}$. By substituting into (35), we have

$$\Pr(\zeta > t) \leq \exp(Cs^2 - st), \text{ for all } |s| \leq \frac{c}{\sigma_1},$$

assuming that the singular value σ_i on the diagonal of Σ are ordered in a descending manner. Minimizing the right-hand side in terms of s , we have

$$\Pr(\zeta > t) \leq \begin{cases} \exp\left(-\frac{t^2}{4C}\right), & \text{if } 0 \leq t \leq \frac{2cC}{\sigma_1} \\ \exp\left(-\frac{ct}{\sigma_1} + \frac{c^2C}{\sigma_1^2}\right) & \text{if } t \geq \frac{2cC}{\sigma_1} \end{cases}, \quad (40)$$

which, because of symmetry, leads to

$$\Pr(|\zeta| > t) \leq \begin{cases} 2 \exp\left(-\frac{t^2}{4C}\right), & \text{if } 0 \leq t \leq \frac{2cC}{\sigma_1} \\ 2 \exp\left(-\frac{ct}{\sigma_1} + \frac{c^2C}{\sigma_1^2}\right) & \text{if } t \geq \frac{2cC}{\sigma_1} \end{cases}. \quad (41)$$

This result is rather similar to ours except that the Hanson-Wright inequality comes with two generic constants c and C . These constants are extremely involved, as shown in the original proof [32]. We need to make all constants precise, and thus maintain our full proof with elementary calculations.

Lemma 4. Let ζ_i , $i = 1, 2, \dots, r$ be i.i.d. copies of the random variable ζ defined in Lemma 3. Then if

$$r \geq 8 \cdot 3^{3/2} \cdot \max\{t^{-5/2}, p^{3/2}t^{-1}\}, \quad (42)$$

we have

$$\Pr \left(\left| \frac{1}{r} \sum_{i=1}^r (\zeta_i^2 - \mathbb{E}[\zeta_i^2]) \right| > t \right) \leq 5r \exp \left(\frac{3}{4} p^{1/2} \right) \exp \left(-\frac{1}{2} r^{1/3} t^{1/3} \right), \quad \text{for } t \in [0, 1]. \quad (43)$$

Remark 2. This lemma essentially deals with the tail bound of a random variable that is of quartic form of a Gaussian. According to the definition, ζ is a quadratic form of Gaussians, and thus is a sub-exponential, but this lemma considers ζ^2 . Quadratic form of sub-exponential vectors are studied in [38]. If we directly employ their results (especially their Corollary 1.6) by setting their $\mathbf{A} = \frac{1}{r} \mathbf{Id}_r \in \mathbb{R}^{r \times r}$, we obtain, for sufficiently large r (made precise in the corollary) that

$$\Pr \left(\left| \frac{1}{r} \sum_{i=1}^r (\zeta_i^2 - \mathbb{E}[\zeta_i^2]) \right| > t \right) \leq C \exp \left(-C' \min \left\{ \left(\frac{r^{1/2}t}{\sqrt{\log r}} \right)^{2/3}, (rt)^{1/3} \right\} \right)$$

where C and C' depend on p . We obtain the same power for r and t as this result, and we make the dependence of the constants on p explicit.

Proof. Let E^t be the event defined as follows:

$$E^t \stackrel{\text{def}}{=} \left\{ \frac{1}{r} \sum_{i=1}^r (\zeta_i^2 - \mathbb{E}[\zeta_i^2]) > t \right\}.$$

Due to the symmetry of $\sum_{i=1}^r \zeta_i^2 - \mathbb{E}[\zeta_i^2]$, the probability in (43) is $2\Pr(E^t)$. We now estimate $\Pr(E^t)$. For any fixed large number M , we define the following event, for $i = 1, 2, \dots, r$:

$$E_i^M \stackrel{\text{def}}{=} \{\zeta_i^2 \leq M\} = \{\zeta_i^2 - 1 \leq M - 1\}.$$

Clearly, we have

$$\Pr(E^t) = \Pr(E^t \cap (\cap_{i=1}^r E_i^M)) + \Pr(E^t \cap (\cap_{i=1}^r E_i^M)^c). \quad (44)$$

We now estimate the two terms.

1. For the first term in (44), we note that

$$\Pr(E^t \cap (\cap_{i=1}^r E_i^M)) = \Pr(E^t \mid (\cap_{i=1}^r E_i^M)) \cdot \Pr((\cap_{i=1}^r E_i^M)) \leq \Pr(E^t \mid (\cap_{i=1}^r E_i^M)). \quad (45)$$

Denoting $X_i \stackrel{\text{def}}{=} \zeta_i^2 - \mathbb{E}[\zeta_i^2]$, and realizing that $\mathbb{E}[\zeta_i^2] = 1$ according to (33) of Lemma 3, then $E_i^M = \{X_i \leq M - 1\}$. Estimating (45) now amounts to controlling the probability of $\sum_{i=1}^r X_i > rt$ assuming that $X_i \leq M - 1$ for all $i = 1, 2, \dots, r$. By applying Bernstein's inequality (72), we have

$$\begin{aligned} \Pr(E^t \mid (\cap_{i=1}^r E_i^M)) &= \Pr \left(\sum_{i=1}^r X_i > rt \mid X_i \leq M - 1, i = 1, 2, \dots, r \right) \\ &\leq \exp \left(-\frac{r^2 t^2 / 2}{\sum_{i=1}^r \mathbb{E}[X_i^2] + (M - 1)rt / 3} \right). \end{aligned}$$

From (34) in Lemma 3, we have $\mathbb{E}[X_i^2] \leq 8$, so that

$$\Pr(E^t \mid (\cap_{i=1}^r E_i^M)) \leq \exp \left(-\frac{3rt^2}{48 + 2(M - 1)t} \right), \quad (46)$$

which gives the upper bound of the first term in (44).

2. For the second term in (44), we note that

$$\Pr\left(E^t \cap (\cap_{i=1}^r E_i^M)^c\right) \leq \Pr\left((\cap_{i=1}^r E_i^M)^c\right) = \Pr\left(\cup_{i=1}^r (E_i^M)^c\right) \leq r \Pr\left((E_i^M)^c\right).$$

By applying (32) from Lemma 3, with $t = \sqrt{M}$, we have

$$\Pr((E_i^M)^c) = \Pr(\zeta_i^2 > M) = \Pr(|\zeta_i| > \sqrt{M}) \leq \begin{cases} 2e^{-\frac{(\sqrt{M}-\sqrt{p})^2}{4\sqrt{p}}} & \text{if } p \leq M \leq 4p \\ 2e^{-\frac{(2\sqrt{M}-3\sqrt{p})}{4}} & \text{if } M \geq 4p \end{cases},$$

and thus

$$\Pr(E^t \cap (\cap_{i=1}^r E_i^M)^c) \leq \begin{cases} 2re^{-\frac{(\sqrt{M}-\sqrt{p})^2}{4\sqrt{p}}} & \text{if } p \leq M \leq 4p \\ 2re^{-\frac{(2\sqrt{M}-3\sqrt{p})}{4}} & \text{if } M \geq 4p \end{cases}. \quad (47)$$

By combining (46) and (47) in (44), we have

$$\Pr(E^t) \leq \exp\left(-\frac{3rt^2}{48 + 2(M-1)t}\right) + \begin{cases} 2re^{-\frac{(\sqrt{M}-\sqrt{p})^2}{4\sqrt{p}}} & \text{if } p \leq M \leq 4p \\ 2re^{-\frac{(2\sqrt{M}-3\sqrt{p})}{4}} & \text{if } M \geq 4p \end{cases}. \quad (48)$$

To find a sharp bound of $\Pr(E^t)$, we choose a suitable value of M . We set

$$M = r^{2/3}t^{2/3}, \quad (49)$$

where r satisfies the lower bound (42). Since $r \geq 8 \cdot 3^{3/2} \cdot p^{3/2}t^{-1}$, we have $r^{2/3} \geq 12pt^{-2/3}$, so that

$$M = r^{2/3}t^{2/3} \geq 12p > 4p, \quad (50)$$

so the second case applies in (48). Since $r \geq 3^{3/2} \cdot 2^3 \cdot t^{-5/2}$, we have $r^{2/3} \geq 12t^{-5/3}$, so that

$$Mt = r^{2/3}t^{5/3} \geq 12,$$

so that, for the denominator of the first term in (48), we have

$$48 + 2(M-1)t = 6Mt + 48 - 2t - 4Mt \leq 6Mt. \quad (51)$$

By using these observations in (48), we have for the value (49) that

$$\Pr(E^t) \leq \exp\left(-\frac{1}{2} \frac{rt}{M}\right) + 2r \exp\left(\frac{3}{4}p^{1/2}\right) \exp\left(-\frac{1}{2}M^{1/2}\right). \quad (52)$$

With M defined as in (49), we see that the two exponential terms involving M in this expression are both equal to $\exp(-r^{1/3}t^{1/3}/2)$. Additionally, since $p \geq 1$ and $r \geq 1$, we have $2r \exp(3p^{1/2}/4) > 4$. Thus, from (52), we obtain

$$\Pr(E^t) \leq (5/2)r \exp\left(\frac{3}{4}p^{1/2}\right) \exp\left(-\frac{1}{2}r^{1/3}t^{1/3}\right). \quad (53)$$

We obtain the result by multiplying the right-hand side by 2, as discussed at the start of the proof. \square

Proposition 1 is a direct consequence of Lemmas 2 and 4.

Proof of Proposition 1. For any $y \in \text{Span}\{F \otimes G\}$, denote $\hat{y} = \frac{y}{\|y\|}$, so that $\|\hat{y}\| = 1$. From Lemma 2, we have

$$\|\hat{S}\hat{y}\|^2 \stackrel{d}{\sim} \frac{1}{r} \sum_{i=1}^r \zeta_i^2, \quad \text{where } \zeta_i \stackrel{\text{def}}{=} \xi_i^\top \Sigma \eta_i,$$

where $\xi_i, \eta_i \in \mathbb{R}^p$ are independent Gaussian vectors drawn from $\mathcal{N}(0, \text{Id}_p)$. We have

$$\frac{\|\hat{S}y\|^2 - \|y\|^2}{\|y\|^2} = \frac{\|\hat{S}\hat{y}\|^2 - \|\hat{y}\|^2}{\|\hat{y}\|^2} = \frac{1}{r} \sum_{i=1}^r \zeta_i^2 - 1.$$

By setting $t = \varepsilon$ in (43) from Lemma 4, we have

$$\Pr \left(\left| \frac{\|\mathbf{S}\mathbf{y}\|^2 - \|\mathbf{y}\|^2}{\|\mathbf{y}\|^2} \right| > \varepsilon \right) = \Pr \left(\left| \frac{1}{r} \sum_{i=1}^r (\zeta_i^2 - 1) \right| > \varepsilon \right) \leq 5r \exp \left(\frac{3}{4} p^{1/2} \right) \exp \left(-\frac{1}{2} r^{1/3} \varepsilon^{1/3} \right),$$

conditioned on $r \geq 8 \cdot 3^{3/2} \cdot \max\{\varepsilon^{-5/2}, p^{3/2} \varepsilon^{-1}\}$, as required. \square

5.2. Proof of Theorem 4. Proposition 1 shows the probability of the sketching matrix \mathbf{S} of the form (24) preserving the norm of a fixed given vector in the range space $\text{Range}(\mathbf{F} \otimes \mathbf{G})$. To show the preservation of norm holds true over the entire column space, we follow the construction of [40]. We construct a γ -net over the unit sphere in $\text{Range}(\mathbf{F} \otimes \mathbf{G})$ and show that for r sufficiently large, with high probability, the angles between any vectors in the net will be preserved with high accuracy. Preservation of angles on the γ -net can be translated to the norm preservation over the entire space.

We show in Lemma 5 that angles can be preserved with the sampling matrix \mathbf{S} of the form (24). In Lemma 7, we calculate the cardinality of the γ -net. The fact that preservation of angle leads to the preservation of norms on the space is justified in Lemma 6. The three results can be combined into a proof for Theorem 4, which we complete at the end of the section.

Lemma 5. *Let V be a collection of vectors in \mathbb{R}^n with cardinality $|V| = f$ and let*

$$\tilde{V} \stackrel{\text{def}}{=} \{\mathbf{u} \pm \mathbf{v} : \mathbf{u}, \mathbf{v} \in V\}.$$

Suppose that a random matrix \mathbf{S} preserved norm on V , in the sense that for each $\tilde{\mathbf{v}} \in \tilde{V}$, with probability at least $1 - \delta$, we have

$$|\|\mathbf{S}\tilde{\mathbf{v}}\|^2 - \|\tilde{\mathbf{v}}\|^2| < \varepsilon \|\tilde{\mathbf{v}}\|^2.$$

Then \mathbf{S} preserves the angle between all elements in V with probability at least $1 - 4f^2\delta$, that is,

$$\Pr(|\langle \mathbf{S}\mathbf{u}, \mathbf{S}\mathbf{v} \rangle - \langle \mathbf{u}, \mathbf{v} \rangle| \leq \varepsilon \|\mathbf{u}\| \|\mathbf{v}\|) > 1 - 4f^2\delta, \quad \text{for all } \mathbf{u}, \mathbf{v} \in V.$$

Proof. Without loss of generality, we assume all vectors in V are unit vectors. Because of the assumptions on \mathbf{S} , we have

$$\Pr \left(|\|\mathbf{S}\tilde{\mathbf{v}}\|^2 - \|\tilde{\mathbf{v}}\|^2| < \varepsilon \|\tilde{\mathbf{v}}\|^2 \text{ for all } \tilde{\mathbf{v}} \in \tilde{V} \right) \leq 1 - f^2\delta. \quad (54)$$

Considering $\mathbf{u}, \mathbf{v} \in V$, we denote $\mathbf{s} \stackrel{\text{def}}{=} \mathbf{u} + \mathbf{v} \in \tilde{V}$ and $\mathbf{t} \stackrel{\text{def}}{=} \mathbf{u} - \mathbf{v} \in \tilde{V}$ and use the parallelogram equality:

$$\langle \mathbf{u}, \mathbf{v} \rangle = \frac{1}{4} (\|\mathbf{s}\|^2 - \|\mathbf{t}\|^2), \quad \langle \mathbf{S}\mathbf{u}, \mathbf{S}\mathbf{v} \rangle = \frac{1}{4} (\|\mathbf{S}\mathbf{s}\|^2 - \|\mathbf{S}\mathbf{t}\|^2),$$

so that

$$\langle \mathbf{S}\mathbf{u}, \mathbf{S}\mathbf{v} \rangle - \langle \mathbf{u}, \mathbf{v} \rangle = \frac{1}{4} (\|\mathbf{S}\mathbf{s}\|^2 - \|\mathbf{s}\|^2 - (\|\mathbf{S}\mathbf{t}\|^2 - \|\mathbf{t}\|^2)).$$

From (54), we have, with probability at least $1 - f^2\delta$, for all $\mathbf{u}, \mathbf{v} \in V$

$$\begin{aligned} |\langle \mathbf{S}\mathbf{u}, \mathbf{S}\mathbf{v} \rangle - \langle \mathbf{u}, \mathbf{v} \rangle| &\leq \frac{1}{4} (|\|\mathbf{S}\mathbf{s}\|^2 - \|\mathbf{s}\|^2| + |\|\mathbf{S}\mathbf{t}\|^2 - \|\mathbf{t}\|^2|) \\ &\leq \frac{\varepsilon}{4} (\|\mathbf{s}\|^2 + \|\mathbf{t}\|^2) = \frac{\varepsilon}{4} (\|\mathbf{u} + \mathbf{v}\|^2 + \|\mathbf{u} - \mathbf{v}\|^2) = \frac{\varepsilon}{4} (2\|\mathbf{u}\|^2 + 2\|\mathbf{v}\|^2) = \varepsilon, \end{aligned}$$

which completes the proof. \square

We now define the γ -net, and show that preservation of angles on this net leads to preservation of norms.

Definition 2. Denote the unit sphere in space $\text{Range}(\mathbf{F} \otimes \mathbf{G})$ by \mathcal{S} , that is,

$$\mathcal{S} \stackrel{\text{def}}{=} \left\{ \mathbf{y} \in \mathbb{R}^{n_1 n_2} : \mathbf{y} = (\mathbf{F} \otimes \mathbf{G})\mathbf{x} \text{ for some } \mathbf{x} \in \mathbb{R}^{p^2} \text{ and } \|\mathbf{y}\| = 1 \right\}. \quad (55)$$

For fixed $\gamma \in (0, 1)$, we call \mathcal{G} a γ -net of \mathcal{S} if \mathcal{G} is a finite subset of \mathcal{S} such that for any $\mathbf{y} \in \mathcal{S}$, there exists $\mathbf{w} \in \mathcal{G}$ such that $\|\mathbf{w} - \mathbf{y}\| \leq \gamma$.

The following lemma was presented in [40, Section 2.1].

Lemma 6. Let \mathcal{S} and \mathcal{G} be as in Definition 2, for some $\gamma \in (0, 1)$. Then preservation of angle on \mathcal{G} leads to the preservation of norm in \mathcal{S} . That is, if

$$|\langle \mathbf{S}\mathbf{w}, \mathbf{S}\mathbf{w}' \rangle - \langle \mathbf{w}, \mathbf{w}' \rangle| \leq \varepsilon, \quad \text{for all } \mathbf{w}, \mathbf{w}' \in \mathcal{G}, \quad (56)$$

then

$$|\|\mathbf{S}\mathbf{y}\|^2 - \|\mathbf{y}\|^2| \leq \frac{\varepsilon}{(1 - \gamma)^2}, \quad \text{for all } \mathbf{y} \in \mathcal{S}.$$

The size of the γ -net can also be controlled, as we now show.

Lemma 7. Let \mathcal{S} the the unit sphere of $\mathbf{F} \otimes \mathbf{G}$, defined in (55). Then for any $\gamma \in (0, 1)$, there exists a γ -net \mathcal{G} of \mathcal{S} such that

$$|\mathcal{G}| \leq \left(1 + \frac{2}{\gamma}\right)^{p^2}.$$

Proof. Notice that \mathcal{S} is isometric to the unit Euclidean sphere \mathcal{S}^{p^2-1} , the result follows directly by applying Corollary 4.2.13 of [37]. \square

Finally, we state the proof of Theorem 4, which is obtained from the lemmas in this section together with Proposition 1.

Proof of Theorem 4. Without loss of generality, it suffices to show \mathbf{S} preserves norm with high accuracy and high probability over the unit sphere in $\text{Range}(\mathbf{F} \otimes \mathbf{G})$, defined by

$$\mathcal{S} \stackrel{\text{def}}{=} \left\{ \mathbf{y} \in \mathbb{R}^{n_1 n_2} : \mathbf{y} = (\mathbf{F} \otimes \mathbf{G})\mathbf{x} \text{ for some } \mathbf{x} \in \mathbb{R}^{p^2} \text{ and } \|\mathbf{y}\| = 1 \right\}.$$

Note from Lemma 7 that for given $\gamma \in (0, 1)$, one can construct a γ -net \mathcal{G} of \mathcal{S} of size $f = (1 + \frac{2}{\gamma})^{p^2}$. Given $\varepsilon_1 \in (0, 1/2)$, then on this \mathcal{G} , according to Proposition 1 and Lemma 5, if we assume

$$r \geq 8 \cdot 3^{3/2} \cdot \max\{\varepsilon_1^{-5/2}, p^{3/2} \varepsilon_1^{-1}\} \quad (57)$$

then with probability at least $1 - \delta_2$ with

$$\delta_2 \leq 20rf^2 \exp\left(\frac{3}{4}p^{1/2}\right) \exp\left(-\frac{1}{2}r^{1/3}\varepsilon_1^{1/3}\right) = 20r \left(1 + \frac{2}{\gamma}\right)^{2p^2} \exp\left(\frac{3}{4}p^{1/2}\right) \exp\left(-\frac{1}{2}r^{1/3}\varepsilon_1^{1/3}\right), \quad (58)$$

we have that \mathbf{S} preserves angles, that is,

$$|\langle \mathbf{S}\mathbf{w}, \mathbf{S}\mathbf{w}' \rangle - \langle \mathbf{w}, \mathbf{w}' \rangle| \leq \varepsilon_1, \quad \text{for all } \mathbf{w}, \mathbf{w}' \in \mathcal{G},$$

According to Lemma 6, \mathbf{S} embeds \mathcal{S} , that is,

$$|\|\mathbf{S}\mathbf{y}\|^2 - \|\mathbf{y}\|^2| \leq \varepsilon, \quad \text{for all } \mathbf{y} \in \mathcal{S}, \quad \text{where } \varepsilon \stackrel{\text{def}}{=} \frac{\varepsilon_1}{(1 - \gamma)^2}.$$

First, we need to convert the condition (57) into one involving ε . We obtain

$$r \geq 8 \cdot 3^{3/2} \cdot \max\{\varepsilon^{-5/2}(1 - \gamma)^{-5}, p^{3/2}\varepsilon^{-1}(1 - \gamma)^{-2}\}. \quad (59)$$

Second, we must alter the lower bound on r to ensure that the right-hand side of (58) is smaller than the given value of δ , that is,

$$\delta \geq 20r \left(1 + \frac{2}{\gamma}\right)^{2p^2} \exp\left(\frac{3}{4}p^{1/2}\right) \exp\left(-\frac{1}{2}r^{1/3}\varepsilon^{1/3}(1-\gamma)^{2/3}\right), \quad (60)$$

or equivalently,

$$\log \delta \geq \log 20 + \log r + 2p^2 \log(1 + 2/\gamma) + \frac{3}{4}p^{1/2} - \frac{1}{2}r^{1/3}\varepsilon^{1/3}(1-\gamma)^{2/3}. \quad (61)$$

Note that for $p \geq 6$ and $\gamma \in (0, 1)$, we have $\log 20 < 3 < .1p^2 \log(1 + 2/\gamma)$ and $.75p^{1/2} < .1p^2 \log(1 + 2/\gamma)$. Thus a sufficient condition for (61) is

$$\log \delta \geq \log r + 2.2p^2 \log(1 + 2/\gamma) - \frac{1}{2}r^{1/3}\varepsilon^{1/3}(1-\gamma)^{2/3}. \quad (62)$$

Denoting

$$\alpha \stackrel{\text{def}}{=} \varepsilon^{1/3}(1-\gamma)^{2/3} \quad \text{and} \quad \beta \stackrel{\text{def}}{=} \frac{1}{3} (2.2p^2 \log(1 + 2/\gamma) + |\log \delta|),$$

we have $\alpha \in (0, 1)$ for any $\varepsilon, \gamma \in (0, 1)$. By using these definitions, we see that (62) is equivalent to

$$\frac{\alpha}{6}r^{1/3} - \log r^{1/3} \geq \beta, \quad (63)$$

for which the combination of the following two conditions is sufficient:

$$\frac{\alpha}{12}r^{1/3} - \log r^{1/3} \geq 0, \quad (64a)$$

$$\frac{\alpha}{12}r^{1/3} \geq \beta. \quad (64b)$$

Condition (64b) can be rewritten to

$$r \geq \frac{12^3\beta^3}{\alpha^3} = \frac{4^3}{\varepsilon(1-\gamma)^2} (2.2p^2 \log(1 + 2/\gamma) + |\log \delta|)^3,$$

for which a sufficient condition is

$$r \geq \frac{8.8^3}{\varepsilon(1-\gamma)^2} \log^3(1 + 2/\gamma) (p^2 + |\log \delta|)^3. \quad (65)$$

The condition (64a) requires $h(r^{1/3}) \geq 0$, where $h(x) \stackrel{\text{def}}{=} \frac{\alpha}{12}x - \log x$. Since

$$h'(x) = \frac{\alpha}{12} - \frac{1}{x} \geq 0,$$

we see that h is an increasing function for $x > 12/\alpha$. By noting that

$$h\left(\frac{12}{\alpha^{5/2}}\right) = \alpha^{-3/2} - \log(12) + \frac{5}{2} \log \alpha \geq 0, \quad \text{for } \alpha \in (0, 0.33),$$

and

$$\frac{12}{\alpha^{5/2}} > \frac{12}{\alpha}, \quad \alpha \in (0, 1),$$

we have for $\alpha \in (0, 0.33)$ that

$$h(r^{1/3}) \geq 0, \quad \text{if } r^{1/3} \geq \frac{12}{\alpha^{5/2}},$$

which leads to

$$r \geq \frac{12^3}{\varepsilon^{5/2}(1-\gamma)^5}. \quad (66)$$

We are free to choose $\gamma \in (0, 1)$ in a way that ensures that $\alpha \in (0, .33)$. In fact, by setting $\gamma = 3/4$, we have

$$\alpha = \varepsilon^{1/3}(1/4)^{2/3} < 0.33, \quad \text{for all } \varepsilon \in (0, 0.5).$$

By combining the conditions (66) and (65), and setting $\gamma = 3/4$, we have

$$r \geq \max \left\{ \frac{\bar{C}_1}{\varepsilon^{5/2}}, \frac{\bar{C}_2}{\varepsilon} (p^2 + |\log \delta|)^3 \right\},$$

with $\bar{C}_2 = 8.8^3 \cdot 4^2 \log^3(11/3) \approx 2.4e4$, and $\bar{C}_1 = 12^3 \cdot 4^5$. \square

We could change the weight in the separation of (63) into (64a) and (64b), one could arrive at different (possibly better) constants \bar{C}_1 and \bar{C}_2 in the final expression. However, our priority is to show dependence of r on ε , δ , and p (and *not* n), and optimization of the constants is less important.

6. NUMERICAL TESTS

This section presents some numerical evidence of the effectiveness of our sketching strategies. We test them on general matrices with the tensor structure and a problem directly from EIT (18). We are mostly concerned of the dependence of accuracy on n , r , and p . The computational complexity is rather straightforward and is omitted from discussion. In both tests, the numerical solutions outperform the theoretical predictions, indicating that there is room for improvement in our bounds for r .

6.1. General matrices with tensor structure. To set up the experiment, we generate two matrices $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_p] \in \mathbb{R}^{n_1 \times p}$ and $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_p] \in \mathbb{R}^{n_2 \times p}$ using:

$$\mathbf{F} = \mathbf{U}_\mathbf{F} \Sigma_\mathbf{F} \mathbf{V}_\mathbf{F}^\top \quad \text{and} \quad \mathbf{G} = \mathbf{U}_\mathbf{G} \Sigma_\mathbf{G} \mathbf{V}_\mathbf{G},$$

where $\mathbf{U}_\mathbf{F} \in \mathbb{R}^{n_1 \times p}$, $\mathbf{U}_\mathbf{G} \in \mathbb{R}^{n_2 \times p}$, $\mathbf{V}_\mathbf{F} \in \mathbb{R}^{p \times p}$, and $\mathbf{V}_\mathbf{G} \in \mathbb{R}^{p \times p}$ are generated by taking the QR-decomposition of random matrices with i.i.d Gaussian entries. The diagonal entries of $\Sigma_\mathbf{F}$ and $\Sigma_\mathbf{G}$ are independently drawn from $\mathcal{N}(1, 0.04)$. Matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ is then defined by setting $\mathbf{a}_j = \mathbf{f}_j \otimes \mathbf{g}_j$, with $n = n_1 n_2$. We further generate the reference solution $\mathbf{x}_\text{ref} \in \mathbb{R}^p$ whose entries are drawn from $\mathcal{N}(1, 0.25)$. The right-hand-side vector $\mathbf{b} \in \mathbb{R}^n$ encodes a small amount of noise; we set

$$\mathbf{b} = \mathbf{A} \mathbf{x}_\text{ref} + 10^{-6} \xi.$$

where each entry of ξ is drawn from $\mathcal{N}(0, 1)$. We compute \mathbf{x}^* using (1).

Three sketching strategies will be considered, the first two cases from (10) and (11), and a third standard strategy that does not take account of the tensor structure in \mathbf{A} .

Case 1: Set $\mathbf{S} = \mathbf{P} \otimes \mathbf{Q}$ (normalized), as defined in (10) with entries in $\mathbf{P} \in \mathbb{R}^{r_1 \times n_1}$ and $\mathbf{Q} \in \mathbb{R}^{r_2 \times n_2}$ drawn i.i.d. from $\mathcal{N}(0, 1)$. Notice here that $r = r_1 r_2$.

Case 2: Set $\mathbf{S}_{i,:} = \mathbf{p}_i^\top \otimes \mathbf{q}_i^\top$ (normalized), as defined in (11), with entries in vectors $\{\mathbf{p}_i\}$ and $\{\mathbf{q}_i\}$ drawn i.i.d. from $\mathcal{N}(0, 1)$ for all $i = 1, \dots, r$.

Random Gaussian: $\mathbf{S} = \mathbf{R} \in \mathbb{R}^{r \times n}$ (normalized), with entries in \mathbf{R} drawn i.i.d. from $\mathcal{N}(0, 1)$.

The random Gaussian choice is not practical in this context, but we include it here as a reference.

For these three choices of \mathbf{S} , we compute the solution \mathbf{x}_s^* of the sketched LS problem (2), and compare the sketching solution with the standard least-squares solution. In particular, we evaluate the following relative error

$$\text{Error} = \frac{f(\mathbf{x}_s^*) - f(\mathbf{x}^*)}{f(\mathbf{x}^*)}, \quad \text{with} \quad f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2. \quad (67)$$

For each strategy, we draw 10 independent samples of \mathbf{S} and compute the median relative error. We discuss how this quantity depends on r and n .

Dependence on r . We set $\varepsilon = 0.5$, $\delta = 10^{-3}$, $p = 10$, and $n_1 = n_2 = 10^2$, and choose the following values for r : 256, 1024, 4096, 16384 and 65536. As shown in Figure 1, the relative error for all three strategies decreases as r increases; all are of the order of r^{-1} . The result suggests Case-2 sketching and the Gaussian reference sketching share almost the same accuracy, while Case 1 is slightly worse.

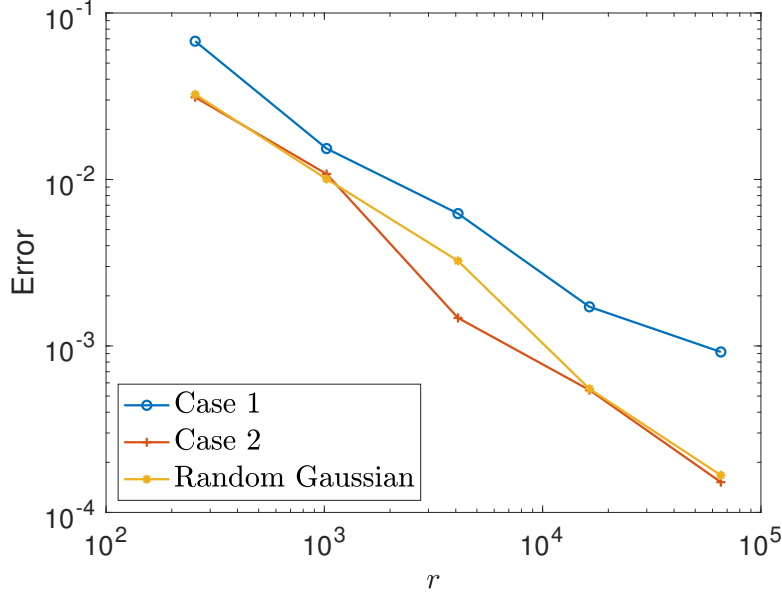


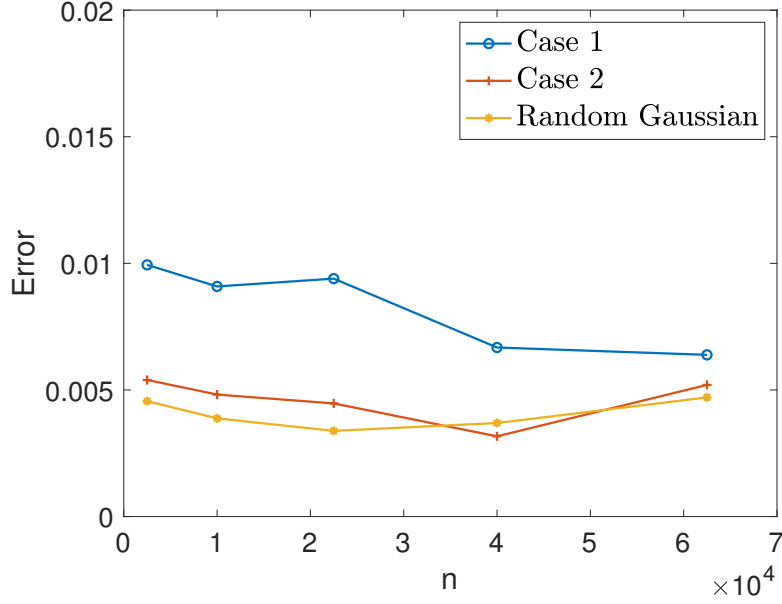
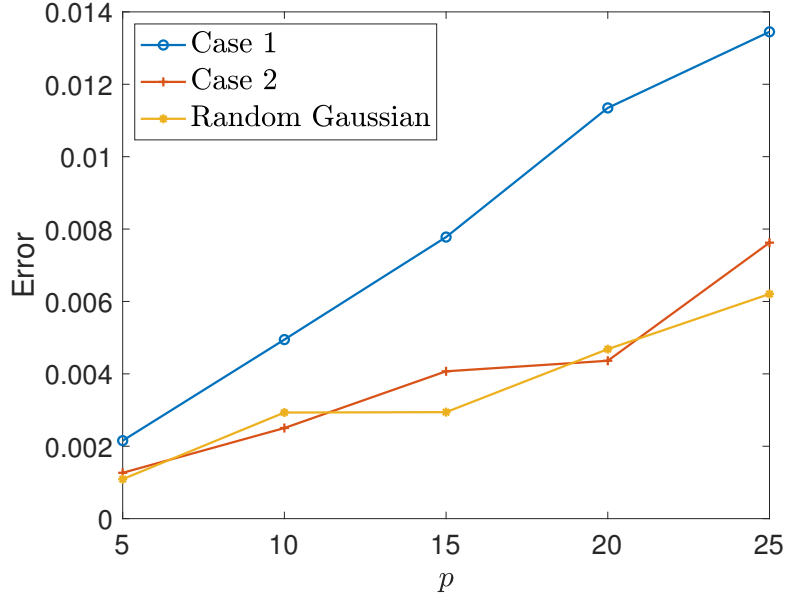
FIGURE 1. Dependence of relative error on r for the three sketching strategies.

Dependence on n . Theorems 3 and 4 suggest essentially no dependence on n . To test this claim empirically, we fix $\varepsilon = 0.5$, $\delta = 10^{-3}$, and $r = 2209$, and set $n_1 = n_2$ to be 50, 100, 150, 200, 250. The error, plotted in Figure 2, shows no dependence on n .

Dependence on p . In this experiment, we study the dependence of relative error on p . We fix $\varepsilon = 0.5$, $\delta = 10^{-3}$, and $r = 4096$ and let p take the values 3, 6, 9, 12, 15. The results are plotted in Figure 3. The plot seems to indicate linear dependence on p , better than the higher powers of p predicted by our bounds. We leave the discussion to future research.

6.2. Electrical Impedance Tomography. In this section, we study the EIT inverse problem on a unit square $[0, 1]^2$. As presented in Section 2, the goal is to reconstruct the conductivity function $\sigma(x)$ in (18). We assume the ground truth $\sigma(x)$ is an indicator function supported at the two yellow squares at the top left and bottom right corners; see Figure 4. The background media $\sigma^*(x)$ (cf. (14)) is set to be a constant function with value 10. We use finite element method to calculate $\rho_1(x)$ and $\rho_2(x)$ on a uniform mesh with $\Delta x = 1/20$. The associated boundary conditions ϕ and ψ are constructed as Dirac-delta functions at all boundary grid points. Under this setup, the matrix A has dimensions $10^4 \times 400$. The right-hand side b is generated by multiplying A with the ground truth $\sigma(x)$ and adding white noise. The EIT inverse problem is highly ill-posed, and thus we set the standard deviation of the mean zero Gaussian noise to be small: 10^{-8} . All three strategies are tested with different number of rows. We record the relative error (67) by taking 10 independent trials.

In Figure 4, we plot the ground truth media $\sigma(x)$ and the reconstructed media using all three different strategies, with $r = 74^2 = 5476$. All of them can roughly reconstruct the unknown function with some

FIGURE 2. Dependence of relative error on ambient dimension n for the three sketching strategies.FIGURE 3. Dependence of relative error on number of unknowns p for the three sketching strategies.

oscillatory errors in the center of the domain. In Figure 5, we plot the relative error in terms of the number of rows r in the sketching matrix S (r is set to be 26^2 , 38^2 , 50^2 , 62^2 , and 74^2). We see that the Case-2 strategy performs as well as the unstructured Gaussian reference, and they both outperform Case 1.

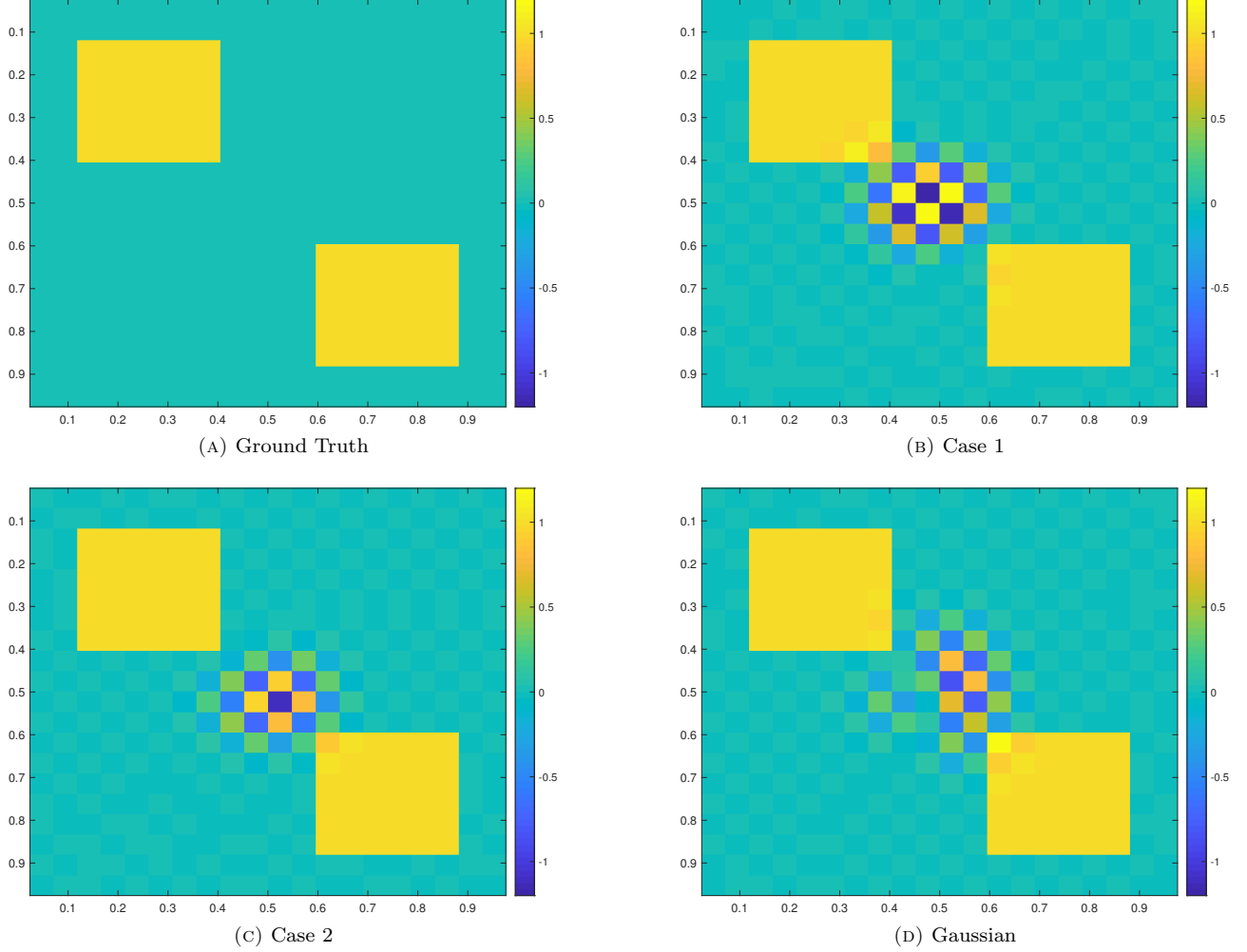


FIGURE 4. The ground truth media and the reconstructed media via all three sketching strategies.

7. CONCLUDING REMARKS

Most PDE-based inverse problems, upon linearization, become Fredholm integral equations, with the testing functions being the product of two functions that are solutions to the forward and the adjoint PDEs. A Khatri-Rao matrix structure arises in the discretization. We study the sketching problem for matrices of this type, where a corresponding structure is enforced in the sketching matrix, for efficiency of computation. We construct the problem under the (ϵ, δ) - l^2 embedding framework, and investigate the number of rows of the sketching matrix that are needed to reconstruct the least-squares solution with ϵ accuracy and δ confidence. The lower bounds differ for the two different sketching strategies that we propose, but both are independent of the size of the ambient space.

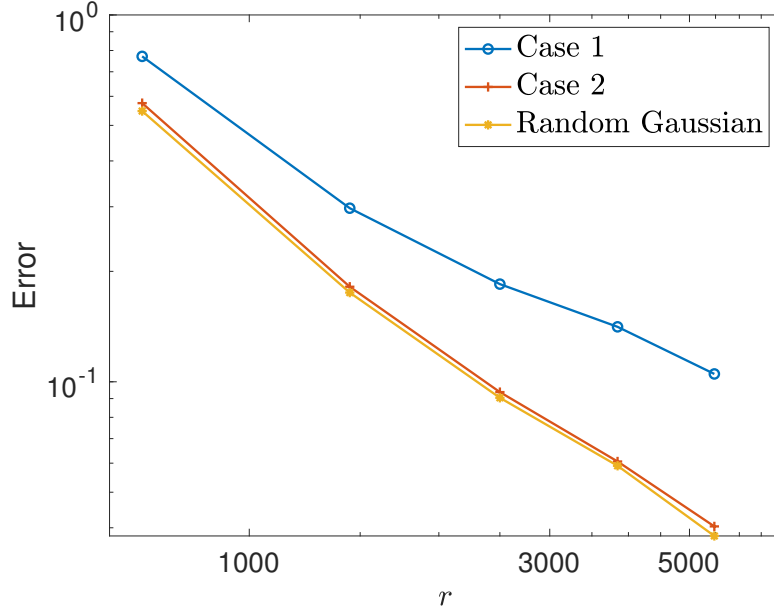


FIGURE 5. For all three strategies, the relative error decreases as the number of rows in \mathbf{S} increases. In particular, the Case-2 sketching strategy performs as well as the unstructured Gaussian strategy.

ACKNOWLEDGMENTS

Chen, Li, and Newton are supported in part by NSF-DMS-1750488 and nsf-tripods 1740707. Wright is supported in part by NSF Awards 1628384, 1634597, and 1740707; Subcontract 8F-30039 from Argonne National Laboratory; and Award N660011824020 from the DARPA Lagrange Program.

APPENDIX A. KEY IDENTITIES AND INEQUALITIES

Some identities and inequalities used repeatedly in the text are collected here.

A.1. Identities of the Kronecker product. Let $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{r_1 \times n_1}$, $\mathbf{B} = (b_{ij}) \in \mathbb{R}^{r_2 \times n_2}$. Then the Kronecker product of \mathbf{A} and \mathbf{B} forms a matrix of size $r_1 r_2 \times n_1 n_2$ defined by:

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n_1}\mathbf{B} \\ a_{21}\mathbf{B} & \cdots & \cdots & a_{2n_1}\mathbf{B} \\ \cdots & \ddots & \ddots & \cdots \\ a_{r_1 1}\mathbf{B} & a_{r_1 2}\mathbf{B} & \cdots & a_{r_1 n_1}\mathbf{B} \end{bmatrix}.$$

The following properties hold.

- (1) Let $\mathbf{A} \in \mathbb{R}^{r_1 \times n_1}$, $\mathbf{B} \in \mathbb{R}^{r_2 \times n_2}$, $\mathbf{C} \in \mathbb{R}^{n_1 \times p_1}$ and $\mathbf{D} \in \mathbb{R}^{n_2 \times p_2}$, then we have the mixed-product property:

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}). \quad (68)$$
- (2) Let $\mathbf{A} \in \mathbb{R}^{r_1 \times n_1}$, $\mathbf{B} \in \mathbb{R}^{r_2 \times n_2}$, and $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$. Further denote by $\mathbf{vec}(\mathbf{X})$ the vectorization of \mathbf{X} formed by stacking the columns of \mathbf{X} into a single column vector, then

$$(\mathbf{B} \otimes \mathbf{A})\mathbf{vec}(\mathbf{X}) = \mathbf{vec}(\mathbf{AXB}^\top). \quad (69)$$

Equivalently, given the same \mathbf{A}, \mathbf{B} and $\mathbf{x} \in \mathbb{R}^{n_1 n_2}$, denote $\mathbf{Mat}(\mathbf{x}) \in \mathbb{R}^{n_1 \times n_2}$ the matricization of the vector \mathbf{x} by aligning subvectors of \mathbf{x} that are of length n_1 into a matrix with n_2 columns, then

$$(\mathbf{B} \otimes \mathbf{A})\mathbf{x} = \mathbf{vec}(\mathbf{A}\mathbf{Mat}(\mathbf{x})\mathbf{B}^\top). \quad (70)$$

A.2. Sub-exponential random variables and Bernstein inequality. Properties of sub-exponential random variables used in the proofs are defined here.

Definition 3. Sub-Exponential random variable A random variable $X \in \mathbb{R}$ is said to be sub-exponential with parameters (λ, b) (denoted as $X \sim \text{subE}(\lambda, b)$) if $\mathbb{E}X = 0$ and its moment generating function satisfies

$$\mathbb{E}e^{sX} \leq \exp\left(\frac{s^2\lambda^2}{2}\right), \quad \text{for all } |s| \leq \frac{1}{b}. \quad (71)$$

We have the following.

Proposition 2. Let $Z \sim \mathcal{N}(0, 1)$, then $X \stackrel{\text{def}}{=} Z^2 - 1$ is sub-exponential with parameters $(2, 4)$.

We conclude with the well known Bernstein inequality.

Proposition 3 (Bernstein inequality). Let X_1, \dots, X_n be i.i.d. mean zero random variables. Suppose that $|X_i| \leq M$ for all $i = 1, \dots, n$, then for any $t > 0$,

$$\Pr\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{t^2/2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + Mt/3}\right). \quad (72)$$

REFERENCES

- [1] S R Arridge, *Optical tomography in medical imaging*, Inverse Problems **15** (1999jan), no. 2, R41–R93.
- [2] Haim Avron, Huy Nguyen, and David Woodruff, *Subspace embeddings for the polynomial kernel*, Advances in neural information processing systems, 2014, pp. 2258–2266.
- [3] Casey Battaglini, Grey Ballard, and Tamara G. Kolda, *A practical randomized cp tensor decomposition*, SIAM Journal on Matrix Analysis and Applications **39** (2018), no. 2, 876–901.
- [4] David J. Biagioni, Daniel Beylkin, and Gregory Beylkin, *Randomized interpolative decomposition of separated representations*, Journal of Computational Physics **281** (2015), 116–134.
- [5] Liliana Borcea, *Electrical impedance tomography*, Inverse Problems **18** (2002oct), no. 6, R99–R136.
- [6] Ke Chen, Qin Li, and Li Wang, *Stability of stationary inverse transport equation in diffusion scaling*, Inverse Problems **34** (2018jan), no. 2, 025004.
- [7] Margaret Cheney, David Isaacson, and Jonathan C. Newell, *Electrical impedance tomography*, SIAM Review **41** (1999), no. 1, 85–101.
- [8] Dehua Cheng, Richard Peng, Yan Liu, and Ioakeim Perros, *Spals: Fast alternating least squares via implicit leverage scores sampling*, Advances in neural information processing systems 29, 2016, pp. 721–729.
- [9] Jocelyn T Chi and Ilse CF Ipsen, *Randomized least squares regression: Combining model- and algorithm-induced uncertainties*, arXiv preprint arXiv:1808.05924 (2018).
- [10] K. Clarkson, P. Drineas, M. Magdon-Ismail, M. Mahoney, X. Meng, and D. Woodruff, *The fast cauchy transform and faster robust linear regression*, SIAM Journal on Computing **45** (2016), no. 3, 763–810.
- [11] Kenneth L Clarkson and David P Woodruff, *Low-rank approximation and regression in input sparsity time*, Journal of the ACM (JACM) **63** (2017), no. 6, 54.
- [12] Huaian Diao, Rajesh Jayaram, Zhao Song, Wen Sun, and David Woodruff, *Optimal sketching for kronecker product regression and low rank approximation*, Advances in neural information processing systems 32, 2019, pp. 4737–4748.
- [13] Huaian Diao, Zhao Song, Wen Sun, and David Woodruff, *Sketching for kronecker product regression and p-splines*, Proceedings of the twenty-first international conference on artificial intelligence and statistics, 201809, pp. 1299–1308.
- [14] Petros Drineas, Michael W. Mahoney, S. Muthukrishnan, and Tamás Sarlós, *Faster least squares approximation*, Numerische Mathematik **117** (2011Feb), no. 2, 219–249.
- [15] Yonina C Eldar and Gitta Kutyniok, *Compressed sensing: theory and applications*, Cambridge university press, 2012.
- [16] Ruhui Jin, Tamara G Kolda, and Rachel Ward, *Faster Johnson-Lindenstrauss transforms via kronecker products*, arXiv preprint arXiv:1909.04801 (2019).

- [17] W. B. Johnson and J. Lindenstrauss, *Extensions of Lipschitz mappings into a Hilbert space*, Contemporary Mathematics **26** (1984), 189–206.
- [18] Michelle Liu, Rajiv Kumar, Eldad Haber, and Aleksandr Aravkin, *Simultaneous shot inversion for nonuniform geometries using fast data interpolation*, arXiv preprint arXiv:1804.08697 (2018).
- [19] Ping Ma, Michael W Mahoney, and Bin Yu, *A statistical perspective on algorithmic leveraging*, The Journal of Machine Learning Research **16** (2015), no. 1, 861–911.
- [20] Michael W. Mahoney, *Randomized algorithms for matrices and data*, Foundations and Trends® in Theoretical Computer Science **3** (2011), no. 2, 123–224.
- [21] Osman Asif Malik and Stephen Becker, *Low-rank tucker decomposition of large tensors using tensorsketch*, Advances in neural information processing systems, 2018, pp. 10096–10106.
- [22] ———, *Guarantees for the kronecker fast Johnson-Lindenstrauss transform using a coherence and sampling argument*, arXiv preprint arXiv:1911.08424 (2019).
- [23] Per-Gunnar Martinsson and Joel Tropp, *Randomized numerical linear algebra: Foundations & algorithms*, arXiv preprint arXiv:2002.01387 (2020).
- [24] Xiangrui Meng and Michael W Mahoney, *Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression*, Proceedings of the forty-fifth annual acm symposium on theory of computing, 2013, pp. 91–100.
- [25] Jelani Nelson and Huy L Nguyễn, *Osnap: Faster numerical linear algebra algorithms via sparser subspace embeddings*, 2013 ieee 54th annual symposium on foundations of computer science, 2013, pp. 117–126.
- [26] Rasmus Pagh, *Compressed matrix multiplication*, ACM Trans. Comput. Theory **5** (August 2013), no. 3, 9:1–9:17.
- [27] M. Pilanci and M. J. Wainwright, *Randomized sketches of convex programs with sharp guarantees*, IEEE Transactions on Information Theory **61** (2015Sep.), no. 9, 5096–5115.
- [28] Mert Pilanci and Martin J. Wainwright, *Iterative hessian sketch: Fast and accurate solution approximation for constrained least-squares*, J. Mach. Learn. Res. **17** (January 2016), no. 1, 1842–1879.
- [29] Garvesh Raskutti and Michael W. Mahoney, *A statistical perspective on randomized sketching for ordinary least-squares*, Journal of Machine Learning Research **17** (2016), no. 213, 1–31.
- [30] Matthew J. Reynolds, Alireza Doostan, and Gregory Beylkin, *Randomized alternating least squares for canonical tensor decompositions: Application to a PDE with random data*, SIAM Journal on Scientific Computing **38** (2016), no. 5, A2634–A2664.
- [31] Vladimir Rokhlin and Mark Tygert, *A fast randomized algorithm for overdetermined linear least-squares regression*, Proceedings of the National Academy of Sciences **105** (2008), no. 36, 13212–13217.
- [32] Mark Rudelson and Roman Vershynin, *Hanson-wright inequality and sub-gaussian concentration*, Electron. Commun. Probab. **18** (2013), 9 pp.
- [33] Tamas Sarlos, *Improved approximation algorithms for large matrices via random projections*, 2006 47th annual ieee symposium on foundations of computer science (focs'06), 2006, pp. 143–152.
- [34] Christian Sohler and David P. Woodruff, *Subspace embeddings for the l_1 -norm with applications*, Proceedings of the forty-third annual acm symposium on theory of computing, 2011, pp. 755–764.
- [35] Yiming Sun, Yang Guo, Joel A Tropp, and Madeleine Udell, *Tensor random projection for low memory dimension reduction*, 2018.
- [36] Roman Vershynin, *Introduction to the non-asymptotic analysis of random matrices*, arXiv preprint arXiv:1011.3027 (2010).
- [37] ———, *High-dimensional probability: An introduction with applications in data science*, Vol. 47, Cambridge university press, 2018.
- [38] Van Vu and Ke Wang, *Random weighted projections, random quadratic forms and random eigenvectors*, Random Structures & Algorithms **47** (2015), no. 4, 792–821.
- [39] David Woodruff and Qin Zhang, *Subspace embeddings and ℓ_p -regression using exponential random variables*, Proceedings of the 26th annual conference on learning theory, 201312, pp. 546–567.
- [40] David P Woodruff, *Sketching as a tool for numerical linear algebra*, Foundations and Trends® in Theoretical Computer Science **10** (2014), no. 1–2, 1–157.