

## Phase retrieval of low-rank matrices by anchored regression

KIRYUNG LEE<sup>†</sup>

*Department of Electrical and Computer Engineering, The Ohio State University, Columbus, OH  
43210, USA*

<sup>†</sup>Corresponding author. Email: lee.8763@osu.edu

SOHAIL BAHMANI

*School of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta,  
GA 30332, USA*

YONINA C. ELGAR

*Department of Computer Science and Applied Mathematics, Weizmann Institute of Science,  
Rehovot 7610001, Israel*

AND

JUSTIN ROMBERG

*School of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta,  
GA 30332, USA*

[Received on 17 December 2018; revised on 29 June 2020; accepted on 3 July 2020]

We study the low-rank phase retrieval problem, where our goal is to recover a  $d_1 \times d_2$  low-rank matrix from a series of phaseless linear measurements. This is a fourth-order inverse problem, as we are trying to recover factors of a matrix that have been observed, indirectly, through some quadratic measurements. We propose a solution to this problem using the recently introduced technique of anchored regression. This approach uses two different types of convex relaxations: we replace the quadratic equality constraints for the phaseless measurements by a search over a polytope and enforce the rank constraint through nuclear norm regularization. The result is a convex program in the space of  $d_1 \times d_2$  matrices. We analyze two specific scenarios. In the first, the target matrix is rank-1, and the observations are structured to correspond to a phaseless blind deconvolution. In the second, the target matrix has general rank, and we observe the magnitudes of the inner products against a series of independent Gaussian random matrices. In each of these problems, we show that anchored regression returns an accurate estimate from a near-optimal number of measurements given that we have access to an anchor matrix of sufficient quality. We also show how to create such an anchor in the phaseless blind deconvolution problem from an optimal number of measurements and present a partial result in this direction for the general rank problem.

### 1. Introduction

We consider the problem of recovering a low-rank matrix  $\mathbf{X}_\#$  from phaseless linear measurements of the form

$$y_m = |\langle \Phi_m, \mathbf{X}_\# \rangle|^2 + \xi_m, \quad m = 1, \dots, M. \quad (1)$$

We refer to this inverse problem as *low-rank phase retrieval* (LRPR). LRPR is a combination of two problems that have received a lot of attention over the past decade. The phase retrieval problem, where

the goal is to recover a vector  $\mathbf{x} \in \mathbb{R}^d$  from  $M$  quadratic measurements of the form  $|\langle \mathbf{x}, \boldsymbol{\phi}_m \rangle|^2$ , is known to be solvable when the  $\boldsymbol{\phi}_m$  are generic and  $M \gtrsim d$  (e.g. see [33] and references therein). There are tractable algorithms for solving the equations that use convex relaxations based on semi-definite programming [14, 16, 61] and polytope constraints [6, 28]. There also exist fast iterative algorithms for nonconvex programming (e.g. [17, 18, 51, 55, 56, 62]). The problem of recovering a  $d_1 \times d_2$  matrix of rank  $r$  from  $M$  linear measurements of the form  $\langle \boldsymbol{\Phi}_m, \mathbf{X} \rangle$  has also been thoroughly analyzed in the literature for generic  $\boldsymbol{\Phi}_m$  [12, 53],  $\boldsymbol{\Phi}_m$  that return samples of the matrix [13, 15, 36, 52] and  $\boldsymbol{\Phi}_m$  with structured randomness [4, 30]; a survey of these results can be found in [22].

Our contribution in this paper is to show that for certain choices of  $\boldsymbol{\Phi}_m$ , we can recover  $\mathbf{X}_\#$  from phaseless measurements (1) from far fewer than  $d_1 d_2$  measurements by taking advantage of the low-rank structure of  $\mathbf{X}_\#$ . Our recovery algorithm uses the recently developed idea of *anchored regression* [6, 7]. The common approaches to estimate  $\mathbf{X}_\#$  from the nonlinear observations (1) lead to nonconvex programs. The anchored regression, however, enables estimation by convex programming as follows. The first step is effectively relaxing the nonlinear equations (1) to convex feasibility constraints. The second step is to use an *anchor matrix*  $\mathbf{X}_0$ , which serves as an initial guess for the solution, to formulate a simple convex program that finds a matrix that is feasible in the relaxed constraints and is best aligned with  $\mathbf{X}_0$ . When the measurements are noiseless ( $\xi_m = 0$ ), we solve

$$\begin{aligned} & \text{minimize}_{\mathbf{X}} \quad -\text{Re} \langle \mathbf{X}_0, \mathbf{X} \rangle + \lambda \|\mathbf{X}\|_* \\ & \text{subject to} \quad |\langle \boldsymbol{\Phi}_m, \mathbf{X} \rangle|^2 \leq y_m, \quad m = 1, \dots, M. \end{aligned} \tag{2}$$

This is a convex program over the space of  $d_1 \times d_2$  matrices. Geometrically, each constraint  $|\langle \boldsymbol{\Phi}_m, \mathbf{X} \rangle|^2 \leq y_m$  is a convex set that has the target  $\mathbf{X}_\#$  on its surface. The program finds an extreme point of the intersection of these convex sets by minimizing the linear functional  $-\text{Re} \langle \mathbf{X}_0, \mathbf{X} \rangle$  regularized by the nuclear norm  $\|\mathbf{X}\|_*$  to account for the low-rank structure of the solution. The success of this program in recovering the target (to within a global phase ambiguity) depends on the behavior of the constraints around  $\mathbf{X}_\#$  and having an anchor  $\mathbf{X}_0$  sufficiently correlated with  $\mathbf{X}_\#$ .

When there is noise, we relax the constraints in (2) and solve

$$\begin{aligned} & \text{minimize}_{\mathbf{X}} \quad -\text{Re} \langle \mathbf{X}_0, \mathbf{X} \rangle + \lambda \|\mathbf{X}\|_* \\ & \text{subject to} \quad \frac{1}{M} \sum_{m=1}^M (|\langle \boldsymbol{\Phi}_m, \mathbf{X} \rangle|^2 - y_m)_+ \leq \eta, \end{aligned} \tag{3}$$

where  $(\cdot)_+$  denotes the positive part function. This yields a stable solution in the sense that if the conditions for noise-free recovery are met, and we choose  $\eta$  larger than the positive part of the perturbations, that is,

$$\eta = \frac{1}{M} \sum_{m=1}^M (-\xi_m)_+ + \epsilon, \quad \text{for some } \epsilon \geq 0,$$

then the solution  $\widehat{\mathbf{X}}$  to (3) obeys  $\|\widehat{\mathbf{X}} - e^{j\theta} \mathbf{X}_\#\|_F \lesssim \eta$  for some  $\theta \in [0, 2\pi)$ . Here  $\epsilon$  denotes an error in estimating the average of the positive part of perturbations by  $\eta$ .

We analyze two scenarios in detail. In the first scenario, the target matrix  $\mathbf{X}_\#$  is of rank  $r$ , and the measurement matrices  $\boldsymbol{\Phi}_m$  have independent real-valued Gaussian entries. Theorem 4.2 below shows that if we start with an anchor matrix that is sufficiently close to  $\mathbf{X}_\#$ , exact recovery occurs when  $M \gtrsim r(d_1 + d_2) \log(d_1 + d_2)$ . Lemma 4.5 shows that the anchor matrix can be computed from the data by a

variation of spectral initialization [17] when the number of measurements  $M$  satisfies  $M \gtrsim r^3 \kappa^4 (d_1 + d_2) \log(d_1 + d_2)$ , where  $\kappa$  denotes the condition number of  $X_{\sharp}$ . We also show that the recovery procedure is stable in the presence of noise.

In our second scenario, the target matrix has rank one,  $X_{\sharp} = \sigma \mathbf{u} \mathbf{v}^*$  with  $\mathbf{u} \in \mathbb{C}^{d_1}$ ,  $\mathbf{v} \in \mathbb{C}^{d_2}$ ,  $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$ , as do the measurement matrices,  $\Phi_m = \mathbf{a}_m \mathbf{b}_m^*$ . As we discuss below, this scenario is a model for the blind deconvolution of two signals from magnitude measurements in the frequency domain. Our analysis in Theorem 4.1 below takes the  $\mathbf{a}_m$  and  $\mathbf{b}_m$  to be complex-valued independent Gaussian random vectors. Under this model, we show that anchored regression produces a stable estimate of  $(\mathbf{u}, \mathbf{v})$  when  $M$  is within a logarithmic factor of  $d_1 + d_2$ . Lemma 4.3 gives a computationally efficient technique for constructing the anchor in a commensurate number of measurements.

## 2. Application: blind deconvolution from Fourier magnitude observations

LRPR arises in a variation of the blind deconvolution, which estimates two unknown signals from the Fourier magnitudes of the convolution. While blind deconvolution is itself an ill-posed, nonlinear problem, the absence of phase information in the Fourier measurements makes it even more challenging. The type of phaseless blind deconvolution problem we describe below arises in various applications in communications and imaging. In optical communications, high spectral efficiency and robustness against adversarial channel conditions for multiple-input multiple-output channels can be achieved using orthogonal frequency division multiplexing (OFDM). Calibrating these communication channels involves solving a blind deconvolution problem. This problem has to be solved from phaseless observations, as practical direct detection receivers work with intensity-only measurements [5] to provide robustness against synchronization errors, which has been one of the key issues in the OFDM systems [10, 54].

A similar calibration problem arises in Fourier ptychography [25]. In this application, an image is computed from phaseless Fourier domain measurements. If there is uncertainty in the point spread function of the optical system, recovering the image becomes a phaseless blind deconvolution problem.

Blind deconvolution that identifies unknown signals  $\mathbf{x}, \mathbf{h} \in \mathbb{C}^M$  (up to reciprocal scaling) from their circular convolution is in general ill-posed but can be solved with a priori information on  $\mathbf{x}$  and  $\mathbf{h}$ . The circular convolution of  $\mathbf{x}$  and  $\mathbf{h}$  can be equivalently expressed in the Fourier domain as the element-wise product, namely

$$\mathbf{F}(\mathbf{x} \circledast \mathbf{h}) = \sqrt{M} \mathbf{F} \mathbf{x} \odot \mathbf{F} \mathbf{h}, \tag{4}$$

where  $\mathbf{F} \in \mathbb{C}^{M \times M}$  is the unitary discrete Fourier matrix of size  $M$ .

We will impose subspace priors on  $\mathbf{x}$  and  $\mathbf{h}$ , modeling  $\mathbf{x} \in \mathbb{C}^M$  as being in the low-dimensional column space of  $\mathbf{D} \in \mathbb{C}^{M \times d_1}$  and  $\mathbf{h}$  as being in the column space of  $\mathbf{E} \in \mathbb{C}^{M \times d_2}$ . Then  $\mathbf{x}$  and  $\mathbf{h}$  are represented as

$$\mathbf{x} = \mathbf{D} \mathbf{u} \quad \text{and} \quad \mathbf{h} = \mathbf{E} \bar{\mathbf{v}}, \tag{5}$$

for some  $\mathbf{u} \in \mathbb{C}^{d_1}$  and  $\mathbf{v} \in \mathbb{C}^{d_2}$ . Here  $\bar{\mathbf{v}}$  denotes the entry-wise complex conjugate of  $\mathbf{v}$ . Let  $\mathbf{a}_m$  denote the  $m$ th column of  $\mathbf{D}^* \mathbf{F}^*$  and  $\mathbf{b}_m$  denote the  $m$ th column of  $\mathbf{E}^{\top} \mathbf{F}^{\top}$  for  $m = 1, \dots, M$ . Then the Fourier measurement of the convolution at frequency  $m$  (after an appropriate normalization) is given as  $\mathbf{a}_m^* \mathbf{u} \mathbf{v}^* \mathbf{b}_m$ . Under this subspace model, it suffices to recover  $\mathbf{u}$  and  $\mathbf{v}$ .

In particular applications, the subspace model for  $\mathbf{h}$  might be introduced as a linear approximation of parametric models via principal component analysis. This technique is used for source localization and channel estimation in underwater acoustics [49, 57]. Some analysis in the context of dimensionality reduction of manifolds is provided in [24, 48].

In the scenario where only noisy Fourier magnitudes of the convolution is observed, the corresponding quadratic measurements are given in the form of

$$y_m = |\mathbf{a}_m^* \mathbf{u} \mathbf{v}^* \mathbf{b}_m|^2 + \xi_m, \quad m = 1, \dots, M,$$

where  $\xi_1, \dots, \xi_M$  denote additive noise. Through the lifting reformulation [3] that substitutes  $\mathbf{u} \mathbf{v}^*$  by a rank-1 matrix  $\mathbf{X}_\#$ , the recovery reduces to a LRPR that estimates the unknown rank-1 matrix  $\mathbf{X}_\#$  from its noisy quadratic measurements:

$$y_m = |\langle \mathbf{a}_m \mathbf{b}_m^*, \mathbf{X}_\# \rangle|^2 + \xi_m, \quad m = 1, \dots, M. \quad (6)$$

This is a particular instance of LRPR and generates quadratic measurements with rank-1 matrices  $\mathbf{a}_1 \mathbf{b}_1^*, \dots, \mathbf{a}_M \mathbf{b}_M^*$ .

In other words, the recovery combines blind deconvolution and phase retrieval; hence, it suffers from the ambiguities in both problems. Similar to phase retrieval, the absence of the phases in the measurements makes the reconstruction a nonconvex problem, even after it has been lifted. By themselves, both phase retrieval and blind deconvolution amount to solving a system of quadratic equations. However, the phaseless blind deconvolution problem (6) is a system of fourth-order equations. Below, we will show that this system can indeed be tractably solved under certain randomness assumptions on the considered subspaces.

### 3. Related work

Recovery of a structured signal from nonlinear measurements has received a significant amount of attention in the past decade, particularly in terms of theoretical analysis of various optimization formulations. A prominent example is the phase retrieval problem, which recovers an unknown signal from quadratic measurements. Unique identification of the solution and performance guarantees of optimization algorithms in the case where the unknown signal is sparse have been recently studied in [7, 11, 19, 26, 34, 39, 45].

Another example, discussed in the previous section, is the blind deconvolution problem, which amounts to solving a system of bilinear equations. Although many approaches for blind deconvolution and its variations have been proposed in the communications, signal processing and computational imaging literature, there has been significant progress in recent years in identifying provable performance guarantees. These results offer theoretical guarantees on the number of measurements  $M$  in (4) as a function of the subspace dimensions  $d_1, d_2$  (number of columns of  $\mathbf{D}, \mathbf{E}$  in (5)) needed to recover  $\mathbf{u}, \mathbf{v}$ . Results that exhibit near-optimal scaling of  $M$  versus  $d_1, d_2$  are known both for convex relaxations of the problem and for iterative algorithms that minimize a nonconvex loss [3, 32, 44]. These results have also been extended to sparsity (in place of subspace) models where the recovery is performed through alternating minimization [42]; however, the near optimal result in this work makes some technical, and perhaps too restrictive, assumptions on the success of projection steps.

The blind deconvolution problem can be made easier if we have the freedom to obtain diversified observations. Specifically, the identification of unknown channel impulse responses excited by an

unknown source has been studied extensively in the communications literature since the 1990s (e.g. [50, 65]). These classical results assumed that the channel responses have finite length and provided algebraic performance guarantees. In recent years, its generalization to the blind gain and phase calibration problem has been analyzed and robust optimization algorithms were proposed along with performance guarantees [2, 21, 41, 43, 46, 47, 63]. There also exists further generalization to the off-the-grid sparsity models [20, 66].

The nonlinear recovery problem considered in this paper is motivated to study a version of blind deconvolution where the convolution measurements are observed through certain nonlinearities. Bendory *et al.* [9] studied a similar problem arising in blind ptychography and identified a set of conditions under which a signal can be identified uniquely from the magnitudes of a short-time Fourier transform taken with an unknown window. In this paper, we are more interested in the recovery by a practical convex program from Fourier magnitudes. The lifting reformulation renders the reconstruction problem into phase retrieval of a low-rank matrix.

The problem of recovering a low-rank matrix from phaseless linear measurements can also be interpreted as a generalization of classical subspace learning (i.e. principal components analysis). This connection was made explicit in [19], where the problem of estimating a covariance matrix from compressed, streaming data was considered. In a subsequent work, [59] considered the quadratic subspace learning problem in a more general setting. A regularized gradient descent method was proposed to solve the LRPR problem, and they provided an analysis for the accuracy of the initialization step under certain randomness assumptions on the measurement matrices.

Unlike the aforementioned works [19, 59], we take a different approach to solving the LRPR problem that uses the recently introduced *anchored regression* [6, 28] technique for relaxing nonlinear measurements. Unlike lifting techniques, this method recasts phase retrieval as a convex program without increasing the number of optimization variables. Unlike techniques based on nonconvex optimization, its analysis relies only on geometry rather than the trajectory of a certain sequence of iterates, which significantly simplifies the derivations. The anchored regression formulation also makes it straightforward to incorporate structural priors on the data through the introduction of convex regularizers [7]. Importantly, we present performance guarantees for stable recovery of low-rank matrices from its random quadratic measurements, which implies exact reconstruction in the noiseless case. Previously, it was only shown that the initialization by a truncated spectral method provides an accurate approximation [59]. After an early version of this paper [40], another approach to the same problem was independently studied in [1]. In contrast to their work, our approach is not restricted to the case of rank-1 measurement matrices and more importantly, anchored regression provides flexibility that allows nonlinearities in the measurement model beyond the quadratic function.

While a general theory for solving equations with convex nonlinearities has been developed, of which (1) is an example, it still remains to compute the key estimates that depend on the structure of the problem (the low-rankness in our case). Furthermore, it is crucial to design an appropriate initialization scheme that provides a valid anchor matrix. We propose a unified approach to the initialization that takes advantage of the separability of the unknown matrix.

It would be of independent interest to see various estimates on functions of random matrices by the noncommutative Rosenthal inequality [35]. All of the matrix Bernstein inequalities [37, 58] and noncommutative Rosenthal inequality [35] provide tail estimates of a sum of independent random matrices. In applying the matrix Bernstein inequalities, one has to verify that all summands have bounded spectral norms (deterministically or almost surely) or compute their Orlicz norms. On the contrary, the noncommutative Rosenthal inequality [35] first computes moment bounds and then provides a tail estimate by the Markov inequality. Particularly when random matrices are given by a

set of Gaussian random variables, the spectral norm is not bounded almost surely and computing the Orlicz norm of the spectral norm is not trivial. Therefore, it is desirable to derive relevant tail estimates by using the noncommutative Rosenthal inequality. Additionally, the expectations of high-order tensor products of Gaussian random vectors in the appendix might be useful in the study of other applications sharing similar tensor structures.

#### 4. Main results

We have four main results. The first two, presented in Section 4.1, give *sample complexity* bounds that relate the accuracy of the estimate returned by (3) to the number of measurements  $M$  observed as in (1). In both cases where the  $\Phi_m$ s are rank-1 matrices such that

$$\Phi_m = \mathbf{a}_m \mathbf{b}_m^*, \quad \mathbf{a}_m \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}), \quad \mathbf{b}_m \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}), \quad m = 1, \dots, M, \quad (7)$$

and when they have i.i.d. standard Gaussian entries, i.e.

$$\text{vec}(\Phi_m) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad m = 1, \dots, M, \quad (8)$$

where  $\text{vec}(\Phi_m)$  denotes the column vector obtained by vertically stacking the columns of  $\Phi$ , we achieve a sample complexity that scales nearly optimally with the size of the target matrix  $X_\#$  and its rank. These results assume that we have an anchor matrix  $X_0$  that is sufficiently correlated with  $X_\#$ .

Our next two main results, presented in Section 4.2, show how such an anchor matrix can be created from the measurements using a spectral initialization method. For the random rank-1 measurements, we are able to construct a sufficiently accurate anchor from a number of observations  $M$  that is proportional to the degrees of freedom in the model of  $X_\#$  up to a logarithmic factor. In the case of Gaussian measurements, we only have a partial result in general and show that a very rough anchor can be bootstrapped into a more accurate one. In a special case where  $X_\#$  is positive semi-definite or is rank-1, however, the results are near-optimal.

##### 4.1 Sample complexity

We begin by presenting theorems that give guarantees on the accuracy of the solution to the convex program (3) in relation to the number of measurements  $M$ . In both of the theorems below, we assume that we have an anchor matrix  $X_0$  that is roughly aligned with the target  $X_\#$ ; we defer the construction of this anchor to Section 4.2.

We start with the case where  $X_\#$  is rank-1, and the measurements are formed by taking the outer product of two random vectors,  $\Phi_m = \mathbf{a}_m \mathbf{b}_m^*$ . As discussed in Section 2 above, this scenario is motivated by problems that involve blind deconvolution from quadratic measurements. Since these applications typically involve the Fourier transform, we formulate our results using complex-valued vectors and matrices.

**THEOREM 4.1** Let  $X_\# = \sigma_\# \mathbf{u}_\# \mathbf{v}_\#^*$  be a complex rank-1 matrix observed as in (6) with  $\Phi_m = \mathbf{a}_m \mathbf{b}_m^*$  for  $m = 1, \dots, M$ , where  $\mathbf{a}_1, \dots, \mathbf{a}_M$  and  $\mathbf{b}_1, \dots, \mathbf{b}_M$  are independent complex Gaussian random vectors as in (7). Suppose that  $X_0 = \mathbf{u}_0 \mathbf{v}_0^*$  with  $\|\mathbf{u}_0\|_2 = \|\mathbf{v}_0\|_2 = 1$  satisfies

$$\inf_{\theta \in [0, 2\pi)} \left\| \mathbf{u}_0 \mathbf{v}_0^* - e^{i\theta} \mathbf{u}_\# \mathbf{v}_\#^* \right\|_{\text{F}} \leq \delta \quad (9)$$

for  $\delta \leq 0.2$ . Then one can set the regularization parameter in (3) such that there exist numerical constants  $C_1, C_2, C_3$  and a constant  $C_\delta$  that depends only on  $\delta$ , for which the following holds.<sup>1</sup> If

$$\frac{M}{\log^2 M} \geq C_\delta(d_1 + d_2), \tag{10}$$

then the solution  $\widehat{\mathbf{X}}$  to (3) satisfies

$$\inf_{\theta \in [0, 2\pi)} \|\widehat{\mathbf{X}} - e^{i\theta} \mathbf{X}_\# \|_{\text{F}} \leq \frac{C_1}{\|\mathbf{X}_\# \|_{\text{F}}} \left( \frac{1}{M} \sum_{m=1}^M |\xi_m| + \epsilon \right) \tag{11}$$

with probability at least  $1 - e^{-C_3 M}$ . Furthermore, the left and right singular vectors  $\widehat{\mathbf{u}}$  and  $\widehat{\mathbf{v}}$  of  $\widehat{\mathbf{X}}$  satisfy

$$\sin \angle(\widehat{\mathbf{u}}, \mathbf{u}_\#) \vee \sin \angle(\widehat{\mathbf{v}}, \mathbf{v}_\#) \leq \frac{C_2}{\|\mathbf{X}_\# \|_{\text{F}}^2} \left( \frac{1}{M} \sum_{m=1}^M |\xi_m| + \epsilon \right). \tag{12}$$

The sufficient number of measurements for stable recovery of  $\mathbf{X}_\#$  (and hence its factors  $\mathbf{u}_\#$  and  $\mathbf{v}_\#$ ) required by (10), scales nearly optimally. That is, the sufficient number of samples is proportional to the degrees of freedom of the unknown rank-1 matrix, i.e.  $d_1 + d_2$ . In Section 4.2 below, we will see that we can also find  $\mathbf{u}_0, \mathbf{v}_0$  that obey (9) from a comparable number of measurements. Combining these results shows that we can recover a  $d_1 \times d_2$  rank-1 matrix from phaseless rank-1 measurements when  $M$  equals to  $d_1 + d_2$  up to a logarithmic factor.

Our second sample complexity result states a performance bound for (3) when the measurements are unstructured Gaussian random matrices and the target is a  $d_1 \times d_2$  matrix of rank  $r$ . This type of measurement model has served as a standard benchmark in the structured recovery literature, and indeed we do obtain a much tighter bound in this case if the target is well conditioned. To ease the derivation, we state the result for real-valued matrices, but it is straightforward to extend it to the complex-valued case at the cost of making the calculations slightly more involved.

**THEOREM 4.2** Let  $\mathbf{X}_\# \in \mathbb{R}^{d_1 \times d_2}$  be of rank  $r$ ,  $\mathbf{X}_\# = \mathbf{U}_\# \boldsymbol{\Sigma}_\# \mathbf{V}_\#^\top$  denote the compact singular value decomposition (SVD) of  $\mathbf{X}_\#$ , and  $\boldsymbol{\Phi}_1, \dots, \boldsymbol{\Phi}_M \in \mathbb{R}^{d_1 \times d_2}$  be Gaussian random matrices as in (8). Suppose that we have an anchor matrix  $\mathbf{X}_0 = \mathbf{U}_0 \mathbf{V}_0^\top$ , where  $\mathbf{U}_0^\top \mathbf{U}_0 = \mathbf{V}_0^\top \mathbf{V}_0 = \mathbf{I}_r$ , that satisfies

$$\min \left( \left\| \mathbf{U}_0 \mathbf{V}_0^\top - \mathbf{U}_\# \mathbf{V}_\#^\top \right\|_{\text{F}}, \left\| \mathbf{U}_0 \mathbf{V}_0^\top + \mathbf{U}_\# \mathbf{V}_\#^\top \right\|_{\text{F}} \right) \leq \delta \left\| \mathbf{U}_\# \mathbf{V}_\#^\top \right\|_{\text{F}} \tag{13}$$

for  $\delta$  that obeys

$$\frac{\delta}{1 - \lambda} \leq 0.45 (2.8 - \kappa), \tag{14}$$

<sup>1</sup> As shown in the proof of Theorem 4.1, given  $\delta$ , one can choose  $\lambda$  explicitly as  $0.9 - \delta$ . For specific methods of constructing the anchor matrix, an appropriate value of  $\delta$  can be determined.

where  $\kappa$  and  $\lambda$  denote the condition number of  $\mathbf{X}_\#$  and the regularization parameter in (3), respectively. Then there exist universal constants  $C_1, C_2$  and a constant  $C_\delta$  that only depends on  $\delta$  for which the following holds. If

$$M \geq C_\delta r(d_1 + d_2) \log(d_1 + d_2), \quad (15)$$

then the solution  $\widehat{\mathbf{X}}$  to (3) satisfies

$$\min \left( \|\widehat{\mathbf{X}} - \mathbf{X}_\#\|_F, \|\widehat{\mathbf{X}} + \mathbf{X}_\#\|_F \right) \leq \frac{C_1}{\|\mathbf{X}_\#\|_F} \left( \frac{1}{M} \sum_{m=1}^M |\xi_m| + \epsilon \right),$$

with probability at least  $1 - e^{-C_2 M}$ .

Although to the authors' knowledge this is the first result of its kind in the literature, and the bound (15) scales in the rank  $r$  and dimensions  $d_1, d_2$  as well as one could hope, we point out a few ways this result could be improved. First, the condition (14) is very restrictive in the sense that it applies only to matrices with a small condition number. Second, constructing  $\mathbf{U}_0 \mathbf{V}_0^\top$  that obeys (13) is non-trivial; as we will see in Section 4.2 below, we will only really be able to do this with confidence when  $\mathbf{X}_\#$  is positive semi-definite or is rank-1.

#### 4.2 Spectral initialization with partial trace

Our main results, presented as Theorems 4.1 and 4.2 above, give bounds on the number of equations  $M$  that are needed to guarantee that the solution to (3) has a certain accuracy. This accuracy depends on the anchor matrix  $\mathbf{X}_0$  being sufficiently close to the unknown matrix  $\mathbf{X}_\#$ . In both cases, we use  $\mathbf{X}_0 = \mathbf{U}_0 \mathbf{V}_0^*$  as an anchor, where  $\mathbf{U}_0 \boldsymbol{\Sigma}_0 \mathbf{V}_0^*$  is the compact SVD of an approximation of  $\mathbf{X}_\#$ ; we are after  $\mathbf{U}_0, \mathbf{V}_0$ , each with orthonormal columns, such that for some  $\delta > 0$  and a unit modulus  $z$ , we have

$$\|\mathbf{U}_0 \mathbf{V}_0^* - z \mathbf{U}_\# \mathbf{V}_\#^*\|_F \leq \delta \|\mathbf{U}_\# \mathbf{V}_\#^*\|_F. \quad (16)$$

In each of the main theorems below, the bounds on  $M$  scale like  $\delta^{-2}$ , and we will achieve the tightest results when we can take  $\delta$  as a constant independent of the matrix dimensions and rank. In this section, we describe a *data-driven* technique for constructing such an anchor matrix.

To understand the challenges in creating the anchor, let us first recall the now well-known spectral initialization for standard phase retrieval for vectors ( $d_2 = 1$  in the formulation above). In this case, we use the observations  $y_m$  to form the  $d_1 \times d_1$  matrix

$$\widehat{\mathbf{R}} = \frac{1}{M} \sum_{m=1}^M y_m \boldsymbol{\Phi}_m \boldsymbol{\Phi}_m^*, \quad (17)$$

and then use the leading eigenvector of  $\widehat{\mathbf{R}}$  as the anchor matrix  $\mathbf{X}_0$ . The idea is that when  $y_m = |\langle \mathbf{X}_0, \boldsymbol{\Phi}_m \rangle|^2$  and the  $\boldsymbol{\Phi}_m$  are random and drawn independent of one another, the expectation of  $\mathbb{E}[y_m \boldsymbol{\Phi}_m \boldsymbol{\Phi}_m^*]$  has a leading eigenvector that is exactly  $\mathbf{X}_0$ , and for  $M$  large enough, the sum in (17)



provides a good approximation to this expectation. In [17], it is shown that (9) holds for constant  $\delta$  when  $M \gtrsim d_1 \log d_1$ .

We might consider using the same initialization when  $X_{\#}$  and the  $\Phi_m$  are matrices. Using a vectorized version of the above, we can form

$$\widehat{\mathbf{R}} = \frac{1}{M} \sum_{m=1}^M y_m \text{vec}(\Phi_m) \text{vec}(\Phi_m)^*,$$

compute the leading eigenvector, then reshape into a  $d_1 \times d_2$  matrix. We are now guaranteed a good anchor when  $M \gtrsim d_1 d_2 \log(d_1 d_2)$ . The problem, though, is that this bound is independent of the rank of  $X_{\#}$ ; we are interested in recovery results that scale as closely as possible to the intrinsic number of degrees of freedom  $r(d_1 + d_2)$  in our matrix model. Simply finding the largest eigenvector of  $\widehat{\mathbf{R}}$  and then re-arranging into a  $d_1 \times d_2$  matrix will not, by itself, result in a matrix that is rank  $r$ , and there is no known algorithm with provable performance guarantees for finding a rank-constrained matrix that is maximally aligned with the column space of  $\widehat{\mathbf{R}}$  (this is a variation on the ‘Sparse principal component analysis (PCA)’ problem).

Our approach for estimating the anchor matrix will be to estimate the row and column spaces of  $X_{\#}$  individually. We will find a  $d_1 \times r$  matrix  $U_0$  whose columns are orthonormal and approximately span the column space, a  $d_2 \times r$  matrix  $V_0$  whose columns are orthonormal and approximately span the row space and then take

$$X_0 = U_0 V_0^*.$$

For the column space estimate  $U_0$ , we choose  $d_2 \times q$  compression matrices  $\Psi_m$  and form

$$\mathbf{Y} = \frac{1}{M} \sum_{m=1}^M y_m \Phi_m \Psi_m \Psi_m^* \Phi_m^*, \tag{18}$$

then take the  $r$  leading eigenvectors of  $\mathbf{Y}$  as  $U_0$ . Similarly for the row space, we choose  $d_1 \times q$   $\Psi'_m$ , form

$$\mathbf{Y}' = \frac{1}{M} \sum_{m=1}^M y_m \Phi_m^* \Psi'_m \Psi'^{*}_m \Phi_m, \tag{19}$$

and take the  $r$  leading eigenvectors as  $V_0$ .

With the measurement matrix  $\Phi_m$  random, we want to choose the compression matrices  $\Psi_m$  in (18) to meet two criteria:

1. The expectation  $\mathbb{E}[\mathbf{Y}]$  has leading eigenvectors that span the same  $r$ -dimensional space as the eigenvectors of  $X_{\#} X_{\#}^*$ .
2. The spectral gap between the  $r$ th and  $(r + 1)$ th eigenvalues of  $\mathbb{E}[\mathbf{Y}]$  is large enough so that it upper bounds the perturbation error  $\|\mathbf{Y} - \mathbb{E} \mathbf{Y}\|$  for relatively small  $M$ . This allows us to use the classical Davis–Kahan theorem to show that the leading eigenvectors of  $\mathbf{Y}$  are approximately aligned with the leading eigenvectors of  $\mathbb{E}[\mathbf{Y}]$ .

Similar statements hold for the  $\Psi'_m$  in (19).

For our blind deconvolution from phaseless measurements application, where  $\mathbf{X}_\# = \sigma \mathbf{u} \mathbf{v}^*$  and  $\Phi_m = \mathbf{a}_m \mathbf{b}_m^*$ , there is a clear way to meet these criteria. If we take

$$\Psi_m = \frac{\mathbf{b}_m}{\|\mathbf{b}_m\|_2^2} \quad \text{and} \quad \Psi'_m = \frac{\mathbf{a}_m}{\|\mathbf{a}_m\|_2^2}, \quad m = 1, \dots, M,$$

then

$$\mathbf{Y} = \frac{1}{M} \sum_{m=1}^M y_m \mathbf{a}_m \mathbf{a}_m^* = \frac{\sigma^2}{M} \sum_{m=1}^M |\mathbf{a}_m^* \mathbf{u}|^2 |\mathbf{v}^* \mathbf{b}_m|^2 \mathbf{a}_m \mathbf{a}_m^* + \xi_m \mathbf{a}_m \mathbf{a}_m^*, \quad (20)$$

$$\mathbf{Y}' = \frac{1}{M} \sum_{m=1}^M y_m \mathbf{b}_m \mathbf{b}_m^* = \frac{\sigma^2}{M} \sum_{m=1}^M |\mathbf{a}_m^* \mathbf{u}|^2 |\mathbf{v}^* \mathbf{b}_m|^2 \mathbf{b}_m \mathbf{b}_m^* + \xi_m \mathbf{b}_m \mathbf{b}_m^*. \quad (21)$$

For independent  $\mathbf{a}_m, \mathbf{b}_m$  that follow (7), a simple calculation yields

$$\mathbb{E} \mathbf{Y} = \sigma^2 \mathbf{u} \mathbf{u}^* + \left( \sigma^2 + \frac{1}{M} \sum_{m=1}^M \xi_m \right) \mathbf{I}, \quad \mathbb{E} \mathbf{Y}' = \sigma^2 \mathbf{v} \mathbf{v}^* + \left( \sigma^2 + \frac{1}{M} \sum_{m=1}^M \xi_m \right) \mathbf{I}.$$

The leading eigenvector for  $\mathbf{Y}$  is the left singular vector  $\mathbf{u}$  for  $\mathbf{X}_\#$ , the leading eigenvector of  $\mathbf{Y}'$  is the right singular vector  $\mathbf{v}$ , and the spectral gap in both cases is  $\sigma^2$ . That  $\mathbf{Y} - \mathbb{E} \mathbf{Y}$  and  $\mathbf{Y}' - \mathbb{E} \mathbf{Y}'$  are small enough so that their leading eigenvectors are close to  $\mathbf{u}$  and  $\mathbf{v}$  when  $M$  is within a logarithmic factor of  $(d_1 + d_2)$  is essentially the content of the following lemma.

**LEMMA 4.3** Let  $\Phi_m = \mathbf{a}_m \mathbf{b}_m^*$  be as in (7). Let  $\mathbf{u}_0 \in \mathbb{C}^{d_1}$  (resp.  $\mathbf{v}_0 \in \mathbb{C}^{d_2}$ ) be the leading eigenvector of  $\mathbf{Y}$  in (20) (resp.  $\mathbf{Y}'$  in (21)) with measurements  $y_m$  constructed as in (1). Let  $\mathbf{u}_\#$  and  $\mathbf{v}_\#$  denote the left and right singular vectors of the rank-1 matrix  $\mathbf{X}_\#$ . Let  $\delta \in (0, 1)$  and  $\alpha \in \mathbb{N}$ . There exist numerical constants  $C_1, C_2$  that only depend on  $\alpha$ , for which the following holds. If

$$\frac{M}{\log^3 M} \geq C_1 \delta^{-2} (d_1 + d_2) \quad (22)$$

and

$$\max_{1 \leq m \leq M} |\xi_m| \leq C_2 \|\mathbf{X}_\#\|^2 \log M, \quad (23)$$

then (9) holds with probability at least  $1 - M^{-\alpha}$ .

**REMARK 4.4** The inequality (23) requires that the signal-to-noise-ratio is larger than the given threshold. The proof of Lemma 4.3 presents a stronger result that holds by (22) and

$$\frac{M}{\log M} \geq C_2 \alpha \left( \frac{\max_{1 \leq m \leq M} |\xi_m|}{\|\mathbf{X}_\#\|^2} \vee \delta^{-1} \left( \frac{\max_{1 \leq m \leq M} |\xi_m|}{\|\mathbf{X}_\#\|^2} \right)^2 \right) \delta^{-1} (d_1 + d_2). \quad (24)$$

Indeed, (23) together with (22) implies (24). Even if (23) is violated, (9) still holds with high probability whenever  $M$  is large enough to satisfy (24) that naturally adapts to the signal-to-noise-ratio. To achieve the order of the logarithmic term in (22), it is necessary to satisfy  $M \lesssim e^{d_1+d_2}$ . Since this upper bound is rather trivial compared to (22), we omit the condition in the statement of Lemma 4.3.

Lemma 4.3 along with Theorem 4.1 give us a clean solution to the phaseless blind deconvolution problem. For generic  $\mathbf{a}_m, \mathbf{b}_m$ , the system

$$y_m = |\langle \mathbf{u}, \mathbf{a}_m \rangle \langle \mathbf{b}_m, \mathbf{v} \rangle|^2 + \text{noise}, \quad m = 1, \dots, M,$$

can be (stably) solved for  $\mathbf{u}, \mathbf{v}$  when  $M$  is within a logarithmic factor of  $d_1 + d_2$ , the total number of unknowns.

For phaseless measurements of a  $d_1 \times d_2$  matrix of rank  $r$ , the story is unfortunately not as clean, even when the  $\Phi_m$  in (1) are i.i.d. Gaussian. The following lemma gives us a partial result on our ability to create a data-driven anchor. It shows that given an estimate of the row space, this estimate can be leveraged into an accurate estimate of the column space.

LEMMA 4.5 Let  $\mathbf{X}_\#$  and  $\Phi_m$ s be as in Theorem 4.2. Let  $\widehat{\mathbf{V}} \in \mathbb{R}^{d_2 \times r}$  satisfy  $\widehat{\mathbf{V}}^\top \widehat{\mathbf{V}} = \mathbf{I}_r$ . Suppose that  $\widehat{\mathbf{V}}$  is given a priori and provides an estimate of the row space of  $\mathbf{X}_\#$  so that

$$\left\| (\mathbf{I}_{d_2} - \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top) \mathbf{V}_\# \mathbf{V}_\#^\top \right\| \leq \delta_{\text{in}} \tag{25}$$

for some  $\delta_{\text{in}} < 1$ . Take  $\mathbf{Y}$  as in (18) with  $\Psi_m = \widehat{\mathbf{V}}$ , and let the columns of  $\mathbf{U}_0$  be the eigenvectors of  $\mathbf{Y}$  corresponding to the  $r$ -largest eigenvalues. Fix  $\delta_{\text{out}} \in (0, 1)$  and  $\alpha \in \mathbb{N}$ . Then there exist numerical constants  $C_1, C_2$  that only depend on  $\alpha$ , for which the following holds. If

$$\frac{M}{\log^3 M} \geq \frac{C_1 \alpha^3 \kappa^4 r^3 d_1}{\delta_{\text{out}}^2 (1 - \delta_{\text{in}})^2} \tag{26}$$

and

$$\frac{\max_{1 \leq m \leq M} |\xi_m|}{\|\mathbf{X}_\#\|} \lesssim \sqrt{r} \log M \wedge \frac{\kappa^2 r^2 \log^2 M}{\delta_{\text{out}} (1 - \delta_{\text{in}})}, \tag{27}$$

then

$$\left\| (\mathbf{I}_{d_1} - \mathbf{U}_0 \mathbf{U}_0^\top) \mathbf{U}_\# \mathbf{U}_\#^\top \right\| \leq \delta_{\text{out}}, \tag{28}$$

holds with probability  $1 - M^{-\alpha}$ , where  $\kappa$  denotes the condition number of  $\mathbf{X}_\#$ .

REMARK 4.6 If the noise is weak enough to satisfy (27), then (26) implies

$$\frac{M}{\log M} \geq C_2 \alpha \left( \frac{\kappa^2 \max_{1 \leq m \leq M} |\xi_m|}{\delta_{\text{out}} (1 - \delta_{\text{in}}) \|\mathbf{X}_\#\|^2} \vee r \left( \frac{\kappa^2 \max_{1 \leq m \leq M} |\xi_m|}{\delta_{\text{out}} (1 - \delta_{\text{in}}) \|\mathbf{X}_\#\|^2} \right)^2 \right) r d_1. \tag{29}$$

Similarly to Remark 4.4, Lemma 4.5 can also be strengthened by substituting (27) by (29). The signal-to-noise-ratio need not be larger than the threshold in (27) whenever  $M$  also satisfies (29). Indeed, this version of Lemma 4.5 is proved in Appendix D.

Lemma 4.5 shows that one obtains an estimate of the column space of accuracy  $\delta_{\text{out}}$  from a given estimate of the row space of accuracy  $\delta_{\text{in}}$ . Here the accuracy is measured by the sine of the principal angle between two subspaces. The number of measurements  $M$  in (26) that guarantees this result increases as one wishes for a more accurate estimate (smaller  $\delta_{\text{out}}$ ) or the input to the initialization method is less accurate (larger  $\delta_{\text{in}}$ ).

Furthermore, it is straightforward to exchange the roles of  $\mathbf{U}_\#$  and  $\mathbf{V}_\#$  above. If we have an estimate  $\widehat{\mathbf{U}}$  of  $\mathbf{U}_\#$ , then we can form  $\mathbf{Y}'$  as in (19) with  $\Phi_m = \widehat{\mathbf{U}}$ , take its leading eigenvectors and have (under analogous conditions as those in the theorem) an accurate estimate of  $\mathbf{V}_\#$ .

The scaling of the number of measurements in (26) has suboptimal dependence on the rank, but its dependence on the side length of the matrix is linear.

Producing an estimate of  $\mathbf{X}_\#$  from matrices  $\mathbf{U}_0$  and  $\mathbf{V}_0$  whose ranges approximate its row and column spaces is itself non-trivial. It involves solving another phase retrieval problem, finding a diagonal  $\Sigma$  so that

$$y_m \approx |\langle \mathbf{U}_0 \Sigma \mathbf{V}_0^\top, \Phi_m \rangle|^2, \quad m = 1, \dots, M.$$

Although it might be possible to control the error propagation from the estimates  $\mathbf{U}_0, \mathbf{V}_0$  to the solution of the problem above, this analysis appears to be extremely complicated.<sup>2</sup> However, there are two specific scenarios where we can upper-bound the error in estimating  $\mathbf{U}_\# \mathbf{V}_\#^\top$  by the subspace estimation errors.

1. Rank-1 case: let  $\sigma_\# \mathbf{u}_\# \mathbf{v}_\#^\top$  be the SVD of  $\mathbf{X}_\#$ . Let  $\phi := \angle(\mathbf{u}_0, \mathbf{u}_\#)$  and  $\psi := \angle(\mathbf{v}_0, \mathbf{v}_\#)$ . Then

$$\begin{aligned} & \left\| \mathbf{u}_0 \mathbf{v}_0^\top - \mathbf{u}_\# \mathbf{v}_\#^\top \right\|_{\text{F}}^2 \wedge \left\| \mathbf{u}_0 \mathbf{v}_0^\top + \mathbf{u}_\# \mathbf{v}_\#^\top \right\|_{\text{F}}^2 = 2 - 2 \cos \phi \cos \psi \leq 2 - 2 \cos^2(\phi \vee \psi) \\ & = 2 \sin^2(\phi \vee \psi) = \left\| (\mathbf{I}_{d_1} - \mathbf{u}_0 \mathbf{u}_0^\top) \mathbf{u}_\# \right\|_2^2 \vee \left\| (\mathbf{I}_{d_2} - \mathbf{v}_0 \mathbf{v}_0^\top) \mathbf{v}_\# \right\|_2^2. \end{aligned}$$

2. Positive semi-definite case: let  $\mathbf{U}_\# \mathbf{A}_\# \mathbf{U}_\#^\top$  be the SVD of  $\mathbf{X}_\#$ . Then

$$\left\| \mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_\# \mathbf{U}_\#^\top \right\|_{\text{F}}^2 = 2r - 2 \left\| \mathbf{U}_0^\top \mathbf{U}_\# \right\|_{\text{F}}^2 = 2 \left\| (\mathbf{I}_{d_1} - \mathbf{U}_0 \mathbf{U}_0^\top) \mathbf{U}_\# \right\|_{\text{F}}^2 \leq 2r \left\| (\mathbf{I}_{d_1} - \mathbf{U}_0 \mathbf{U}_0^\top) \mathbf{U}_\# \right\|_{\text{F}}^2.$$

For the above two cases, one can combine Theorem 4.2 and Lemma 4.5 to get a complete analysis of the regularized anchored regression. In the latter case, we still assume that an estimate of  $\mathbf{U}_\#$  is given a priori. Lemma 4.5 provides a refined estimate so that we can invoke Theorem 4.2 with the resulting  $\mathbf{U}_0$ .

<sup>2</sup> An alternative approach is to estimate  $\Sigma$  from  $\mathbf{U}_0$  and  $\mathbf{V}_0$  through extra independent random measurements. However, this approach doubles the number of observations and may not be interesting in practice. Therefore, we pursue analysis in some special cases without extra observations.

### 5. Proof of main results

The convex program for phase retrieval of low-rank matrices in (3) is variation to a special case of the anchored regression studied in [7] and the performance guarantees in this paper primarily follow from the main results in [7]. The theorems stated in the previous section are basically obtained by computing the key quantities that determine the sample complexity.

#### 5.1 Theoretical analysis of regularized anchored regression

At the core of our analysis is an adaptation of the main result of [7]. The main idea of [7, Theorem 2.1] is to use the small-ball method to find a uniform lower bound for a certain empirical process that is determined by the independent random matrices  $\Phi_1, \dots, \Phi_M$  and indexed by a *deterministic* set  $\mathcal{H} \subset \mathbb{C}^{d_1 \times d_2}$  containing  $\Delta = \widehat{\mathbf{X}} - \mathbf{X}_\#$ . Then this uniform lower bound implies an upper bound for the estimation error  $\Delta$ .

However, the original statement of [7, Theorem 2.1] cannot be applied directly to the problem of interest in this paper because of two important differences. First, due to technical challenges in our specific problem, as elaborated in Section 4.2, we can only construct a weaker form of anchor compared to that considered originally in [7]. Second, we want to address the case of recovering complex and rank-1 matrices as considered in Theorem 4.1. The results of [7], however, only consider variables and operations in the real space. Therefore, we need to adapt the result of [7] with slight modifications so that it becomes compatible with our setting.

As discussed in Section 4.2, instead of an anchor that approximates the ground truth  $\mathbf{X}_\#$ , we require the anchor to approximate  $\mathbf{U}_\# \mathbf{V}_\#^*$  up to a global phase. To be explicit, we only need to consider a complex phase ambiguity in the case of recovering a complex rank-1 target, where we have  $\mathbf{U}_\# = \mathbf{u}_\#$  and  $\mathbf{V}_\# = \mathbf{v}_\#$  and the anchor should basically approximate  $\mathbf{u}_\# \mathbf{v}_\#^*$ . In the case of recovering a real-valued low-rank matrix, the phase ambiguity simply reduces to a sign ambiguity.

With these consideration in mind, here and throughout, we assume that the global phase of the anchor  $\mathbf{X}_0$  is aligned with  $\mathbf{U}_\# \mathbf{V}_\#^*$ , namely

$$\text{Re} \langle \mathbf{X}_0, \mathbf{U}_\# \mathbf{V}_\#^* \rangle \geq 0, \quad \text{Im} \langle \mathbf{X}_0, \mathbf{U}_\# \mathbf{V}_\#^* \rangle = 0, \tag{30}$$

which, if we operate entirely in the real domain, simply reduces to  $\langle \mathbf{X}_0, \mathbf{U}_\# \mathbf{V}_\#^\top \rangle \geq 0$ . The assumption (30) can be made without loss of generality because of the following *equivariance* property. For any  $\theta \in [0, 2\pi)$ , if we replace the anchor  $\mathbf{X}_0$  in (3) by  $e^{i\theta} \mathbf{X}_0$ , then the original solution  $\widehat{\mathbf{X}}$  accordingly changes to  $e^{i\theta} \widehat{\mathbf{X}}$ . This property is due to fact that the nuclear norm as well as the constraints in (3) are invariant under the mapping  $\mathbf{X} \mapsto e^{i\theta} \mathbf{X}$ . Since we define the accuracy as the distance to the orbit of  $\mathbf{X}_\#$ , i.e.  $\{e^{i\omega} \mathbf{X}_\# : \omega \in [0, 2\pi)\}$ , the mentioned adjustment of the anchor will not affect the accuracy guarantees. Indeed, under (30), the assumption in (16) simplifies to

$$\|\mathbf{X}_0 - \mathbf{U}_\# \mathbf{V}_\#^*\|_F \leq \delta \|\mathbf{U}_\# \mathbf{V}_\#^*\|_F = \delta \sqrt{r}. \tag{31}$$

Since  $\widehat{\mathbf{X}}$  is a minimizer to (3) and  $\mathbf{X}_\#$  is within its feasible set, it naturally follows that  $\Delta = \widehat{\mathbf{X}} - \mathbf{X}_\#$  belongs to the set of all ascent directions of the objective function given by

$$\mathcal{A} := \left\{ \mathbf{H} \in \mathbb{C}^{d_1 \times d_2} : \inf_{\mathbf{G} \in \lambda \partial \|\mathbf{X}_\#\|_*} \text{Re} \langle \mathbf{X}_0 - \mathbf{G}, \mathbf{H} \rangle \geq 0 \right\}.$$

It is desirable to construct the anchor matrix  $\mathbf{X}_0$  from the available measurements and avoid *sample splitting* schemes. However, for such constructions of the anchor matrix, the set  $\mathcal{A}$  will also depend on the measurement matrices  $\{\Phi_m\}_{m=1}^M$  that complicates the analysis. To avoid these complications, similar to the approach of [7], we relax  $\mathcal{A}$  to some superset that is not dependent on the measurement matrices. Here we consider the superset  $\mathcal{A}_\delta$  of  $\mathcal{A}$ , defined as

$$\mathcal{A}_\delta := \left\{ \mathbf{H} \in \mathbb{C}^{d_1 \times d_2} : \inf_{\mathbf{G} \in \lambda \partial \|\mathbf{X}_\sharp\|_*} \sqrt{r} \delta \|\mathbf{H}\|_F + \operatorname{Re} \langle \mathbf{U}_\sharp \mathbf{V}_\sharp^* - \mathbf{G}, \mathbf{H} \rangle \geq 0 \right\}, \quad (32)$$

which is clearly independent of  $\{\Phi_m\}_{m=1}^M$ . Inclusion of  $\mathcal{A}$  in  $\mathcal{A}_\delta$  follows from (31), the triangle inequality and the Cauchy–Schwarz inequality.

To address a technical challenge that only arises when operating in the complex domain, for recovery of complex rank-1 matrices we need to make another modification compared to the original result of [7]. Specifically, similar to [6], with  $\mathbf{X}_\sharp = \sigma_\sharp \mathbf{u}_\sharp \mathbf{v}_\sharp^*$  as the complex rank-1 ground truth, we introduce the set

$$\mathcal{R}_\delta := \left\{ \mathbf{H} \in \mathbb{C}^{d_1 \times d_2} : \left\| \mathbf{H} - \frac{\mathbf{X}_\sharp \langle \mathbf{X}_\sharp, \mathbf{H} \rangle}{\|\mathbf{X}_\sharp\|_F^2} \right\|_F \geq \frac{\sqrt{1 - \delta^2} |\operatorname{Im} \langle \mathbf{X}_\sharp, \mathbf{H} \rangle|}{\delta \|\mathbf{X}_\sharp\|_F} \right\}. \quad (33)$$

Obviously,  $\mathcal{R}_\delta$  is only important if we operate in the complex domain; in the real domain,  $\mathcal{R}_\delta$  is the entire space and effectively can be ignored. The following lemma, proved in Appendix E, the set  $\mathcal{R}_\delta$  also contains  $\mathbf{A}$  when  $\mathbf{X}_0$  and  $\mathbf{X}_\sharp$  are at most  $\delta$ -apart.

LEMMA 5.1 With  $\mathbf{X}_\sharp = \sigma_\sharp \mathbf{u}_\sharp \mathbf{v}_\sharp^*$ , suppose that (30) and

$$\left\| \mathbf{X}_0 - \frac{\mathbf{X}_\sharp \langle \mathbf{X}_\sharp, \mathbf{X}_0 \rangle}{\|\mathbf{X}_\sharp\|_F^2} \right\|_F \leq \delta \|\mathbf{X}_0\|_F \quad (34)$$

hold. Then  $\widehat{\mathbf{X}} - \mathbf{X}_\sharp \in \mathcal{R}_\delta$ .

Finally, based on the arguments in [7, Theorem 2.1], our result depends on the following two key quantities defined with respect to the set  $\mathcal{H} = \mathcal{A}_\delta \cap \mathcal{R}_\delta$ . First, the Rademacher complexity of  $\mathcal{H}$  is defined as

$$\mathfrak{C}_M(\mathcal{H}) := \mathbb{E} \sup_{\mathbf{H} \in \mathcal{H} \setminus \{0\}} \frac{1}{\sqrt{M}} \sum_{m=1}^M \frac{\epsilon_m \operatorname{Re}(\langle \mathbf{X}_\sharp, \Phi_m \rangle \langle \Phi_m, \mathbf{H} \rangle)}{\|\mathbf{H}\|_F}, \quad (35)$$

where  $\epsilon_1, \dots, \epsilon_M$  are i.i.d. Rademacher random variables independent of everything else. Second, for  $\tau > 0$ , we also consider a variation of small-ball probability that is defined as

$$P_\tau(\mathcal{H}) := \inf_{\mathbf{H} \in \mathcal{H}} \mathbb{P}(\operatorname{Re}(\langle \mathbf{X}_\sharp, \Phi_m \rangle \langle \Phi_m, \mathbf{H} \rangle) \geq \tau \|\mathbf{H}\|_F). \quad (36)$$

Equipped with these notions, the following theorem provides the accuracy guarantees for the regularized anchored regression in the context of LRPR problem.

**THEOREM 5.2** (An adaptation of [7, Theorem 2.1 for LRPR]). Suppose that  $\Phi_1, \dots, \Phi_M$  in (3) are independent random matrices, and  $X_0$  satisfies (16), (30), and (34), where  $0 < \delta < 1$ . Recalling the definitions (32), (33), (35) and (36), for any  $t > 0$ , if

$$M \geq 4 \left( \frac{\mathfrak{C}_M(\mathcal{A}_\delta \cap \mathcal{R}_\delta) + t\tau}{\tau P_\tau(\mathcal{A}_\delta \cap \mathcal{R}_\delta)} \right)^2, \tag{37}$$

then the solution  $\widehat{X}$  to (3) obeys

$$\inf_{\theta \in (0, 2\pi)} \|\widehat{X} - e^{i\theta} X_\# \|_F \leq \frac{2}{\tau P_\tau(\mathcal{A}_\delta \cap \mathcal{R}_\delta)} \left( \frac{1}{M} \sum_{m=1}^M |\xi_m| + \epsilon \right)$$

with probability at least  $1 - e^{-2t^2}$ .

**REMARK 5.3** A few remarks on Theorem 5.2 are in order.

1. We emphasize again that the required conditions in (30) for the anchor can be made without loss of generality due to the equivariance property discussed above.
2. The original result in [7, Theorem 2.1] considered the problem only in the real domain, where the condition (30) is reduced to the (implicit) assumption  $\langle X_0, X_\# \rangle \geq 0$ . As mentioned above, in this scenario the set  $\mathcal{R}_\delta$  becomes trivial (i.e.  $\mathcal{R}_\delta = \mathbb{R}^{d_1 \times d_2}$ ) as well.
3. The additive noise  $\xi_m$  to the quadratic measurement  $|\langle \Phi_m, X_\# \rangle|^2$  is arbitrary fixed. Specifically, we assume that  $\xi_m$  does not depend on  $\Phi_1, \dots, \Phi_M$ .

Theorems 4.1 and 4.2 are then obtained from Theorem 5.2 by specifying key estimates depending on the corresponding measurement matrices. For the convenience in computing these estimates, we provide a more explicit characterization of  $\mathcal{A}_\delta$  as follows. The subdifferential of  $\|\cdot\|_*$  at  $X_\#$ , whose SVD is  $U_\# \Sigma_\# V_\#^*$ , is expressed as

$$\partial \|X_\# \|_* = \left\{ Z : \mathcal{P}_T(Z) = U_\# V_\#^*, \|\mathcal{P}_{T^\perp}(Z)\| \leq 1 \right\}, \tag{38}$$

where  $\mathcal{P}_T : \mathbb{C}^{d_1 \times d_2} \rightarrow \mathbb{C}^{d_1 \times d_2}$  denotes the orthogonal projection onto the *tangent space*  $T$  of the manifold of rank- $r$  matrices at  $X_\#$  given by

$$T = \left\{ U_\# \tilde{V}^* + \tilde{U} V_\#^* : \tilde{V} \in \mathbb{C}^{d_2 \times r}, \tilde{U} \in \mathbb{C}^{d_1 \times r} \right\}$$

and  $\mathcal{P}_{T^\perp} : \mathbb{C}^{d_1 \times d_2} \rightarrow \mathbb{C}^{d_1 \times d_2}$  denotes the projection onto  $T^\perp$ , the perpendicular subspace of  $T$ . By plugging in the expression of the subdifferential in (38) to (32), we obtain an alternative expression of  $\mathcal{A}_\delta$  given by

$$\mathcal{A}_\delta = \left\{ H \in \mathbb{C}^{d_1 \times d_2} : \sqrt{r}\delta \|H\|_F - \lambda \|\mathcal{P}_{T^\perp}(H)\|_* + (1 - \lambda) \text{Re} \langle U_\# V_\#^*, H \rangle \geq 0 \right\}. \tag{39}$$

### 5.2 Proof of Theorem 4.2

All matrices and scalars are real-valued in Theorem 4.2. Thus,  $\mathcal{R}_\delta$  becomes trivial and it suffices to compute estimates of  $P_\tau(\mathcal{H})$  and  $\mathfrak{C}_M(\mathcal{H})$  for  $\mathcal{H} = \mathcal{A}_\delta$ . The following lemmas respectively provide estimates of  $P_\tau(\mathcal{A}_\delta)$  and  $\mathfrak{C}_M(\mathcal{A}_\delta)$  whose proofs are deferred to Appendices F and G.

LEMMA 5.4 Suppose the hypotheses in Theorem 4.2 hold. Then for any  $\tau' > 0$ ,

$$\inf_{\mathbf{H} \in \mathcal{A}_\delta} \mathbb{P} \left( \operatorname{Re}(\langle \mathbf{X}_\#^*, \boldsymbol{\Phi}_m \rangle \langle \boldsymbol{\Phi}_m, \mathbf{H} \rangle) \geq \tau' \|\mathbf{X}_\#^*\|_{\mathbb{F}} \|\mathbf{H}\|_{\mathbb{F}} \right) \geq \frac{\exp(-20\tau')}{10}. \quad (40)$$

LEMMA 5.5 Suppose the hypotheses in Theorem 4.2 hold. Then

$$\mathfrak{C}_M(\mathcal{A}_\delta) \leq \frac{C(1 - \lambda + \delta) \|\mathbf{X}_\#^*\|_{\mathbb{F}} \sqrt{r(d_1 + d_2) \log(d_1 + d_2)}}{\lambda} \quad (41)$$

for a numerical constant  $C$ .

To prove Theorem 4.2, we only need to apply the above estimates in Theorem 5.2. We first show that the assumptions of Theorem 4.2 are sufficient to invoke Theorem 5.2. Following the discussion in Section 5.1, the condition (30) can be satisfied without loss of generality by flipping the sign of  $\mathbf{X}_0$  if necessary. Fix  $\tau'$  to a positive constant (e.g.  $\tau' = 0.1$ ). Let  $\tau = \tau' \|\mathbf{X}_\#^*\|_{\mathbb{F}}$ . Then Lemma 5.4 implies that  $\tau P_\tau(\mathcal{A}_\delta) \geq c \|\mathbf{X}_\#^*\|_{\mathbb{F}}$  for a numerical constant  $c > 0$ . Choosing  $\lambda = 0.9 - \delta$  makes the right-hand side of (41) an increasing function of  $\delta$ . Then, by Lemma 5.5, the Rademacher complexity  $\mathfrak{C}_M(\mathcal{A}_\delta)$  is upper-bounded by  $\sqrt{r(d_1 + d_2) \log(d_1 + d_2)}$  up to a constant solely determined by  $\delta$ . Therefore, (15) implies that (37) holds whenever  $t\tau'$  is dominated by  $\sqrt{M}$ . We can choose  $t$  so that the probability of failure is at most  $e^{-2t^2} = e^{-cM}$ , for some numerical constant  $c > 0$ .

### 5.3 Proof of Theorem 4.1

Theorem 4.1 considers recovery of complex-valued rank-1 matrices. We apply Theorem 5.2 for  $\mathcal{H} = \mathcal{A}_\delta \cap \mathcal{R}_\delta$  to prove Theorem 4.1. The following lemmas, proved in Appendices H and I, respectively, provide a lower bound on  $P_\tau(\mathcal{H})$  and an upper bound on  $\mathfrak{C}_M(\mathcal{H})$ .

LEMMA 5.6 Suppose the hypotheses in Theorem 4.1 hold. Suppose that  $\delta + \lambda < 1$  and  $\delta \leq 0.2$ . Then there exists a numerical constant  $\tau' > 0$  such that

$$\inf_{\mathbf{H} \in \mathcal{A}_\delta \cap \mathcal{R}_\delta} \mathbb{P} \left( \operatorname{Re}(\mathbf{b}^* \mathbf{X}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{H} \mathbf{b}) \geq \tau' \|\mathbf{X}_\#^*\|_{\mathbb{F}} \|\mathbf{H}\|_{\mathbb{F}} \right) \geq C_{\tau'},$$

where  $C_{\tau'}$  is a positive numerical constant that only depends on  $\tau'$ .

LEMMA 5.7 Suppose the hypotheses in Theorem 4.1 hold. Then

$$\mathfrak{C}_M(\mathcal{A}_\delta) \leq \frac{C(1 - \lambda + \delta) \|\mathbf{X}_\#^*\|_{\mathbb{F}} \sqrt{d_1 + d_2} \log M}{\lambda} \quad (42)$$

for a numerical constant  $C$ .



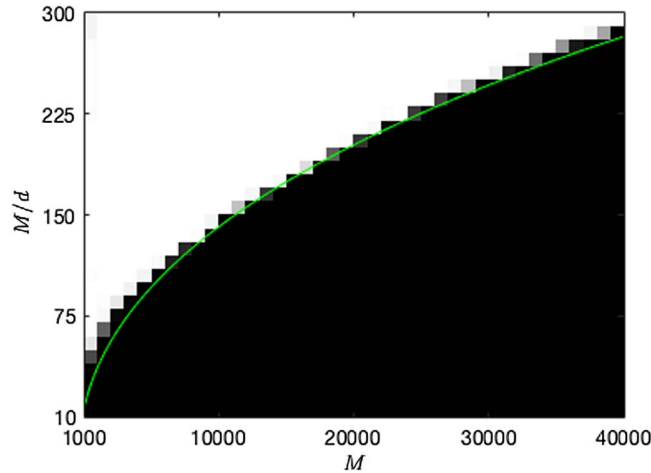


FIG. 1. Empirical phase transition in the noiseless case with rank-1 measurements. The success rate out of 100 trials is plotted in a gray scale (white: all success, black: all failure).

The error bound in (11) then follows from Theorem 5.2 with the above estimates given by Lemmas 5.6 and 5.7. To apply Lemma 5.6, we choose  $\lambda = 0.9 - \delta$ . Then similar to the proof of Theorem 4.2, the factor  $(1 - \lambda + \delta)/\lambda$  becomes an increasing function in  $\delta$ . The constant  $C_\delta$  is given by this function of  $\delta$  together with the result of Lemma 5.6.

Finally, the error bound for the estimation of  $\mathbf{u}$  and  $\mathbf{v}$  in (12) follows immediately from the Davis–Kahan theorem (Theorem C.1).

### 6. Numerical results

We conducted a Monte Carlo simulation to study the empirical performance of the proposed convex programs. Specifically, we considered the optimization problem in (2) in the noiseless case where the measurement matrices are given as the outer product of two Gaussian random vectors and the unknown rank-1 matrix is a square matrix ( $d_1 = d_2 = d$ ). To solve (2), we used the software package TFOCS [8] that uses a smoothed conic dual formulation.

Figure 1 illustrates the empirical phase transition. For a fixed number of measurements  $M$ , we vary the matrix size  $d$  where the ratio  $M/d$  belongs to a given interval. In Fig. 1, the convex program provides the exact recovery when  $d$  is below a certain threshold determined by  $M$ . The sample complexity result by Theorem 4.1 and Lemma 4.3 quantifies this threshold as  $CM/\log^\alpha M$  for some constants  $C, \alpha > 0$ . Alternatively, if the oversampling rate  $M/d$  exceeds a polylog factor of  $M$ , then the convex program provides the exact recovery. The empirical phase transition occurs at  $M/d \approx 0.14 \log^5 M$  or  $d \approx 7.3M/\log^5 M$  indicated by the green curve in the figure. Although the requirements for the constants  $C$  and  $\alpha$  in our proofs seem conservative, our theory is consistent with the empirical performance up to the choice of these constants.

## 7. Discussions

We proposed a simple initial estimation using partial traces. The regularized anchored regression with the nuclear norm given by this initial estimate provides a stable estimate for LRPR. Performance guarantees were derived for several random measurement models.

One may consider several alternative approaches to LRPR. By applying the lifting trick twice, LRPR can be equivalently cast as a linear inverse problem where the solution is a 4-way tensor of rank-1. Then recovery can be carried out by tensor nuclear norm minimization as a convex relaxation. This fully convexified approach can be attractive as it does not require any initial estimate. However, it is NP-hard to compute the tensor nuclear norm [31], which renders the approach impractical. Recent iterative nonconvex methods such as Wirtinger flow [17] and its variation for the case of sparse signals [11] can be adapted to address low-rank recovery problems. Nevertheless, it will be still required to obtain an accurate initial estimate. Furthermore, the analysis might involve more complications such as the need for resampling in each iteration. In contrast, anchored regression, originally proposed for ordinary phase retrieval and later modified to a regularized version to accommodate various priors on the solution, allows a streamlined analysis by leveraging the geometry of the corresponding convex optimization problem. Its low computational cost also competes with that of the nonconvex approaches.

## Acknowledgements

The authors thank the anonymous reviewers for their careful reading of the manuscript and their many insightful comments and suggestions.

## Funding

National Science Foundation Computing and Communication Foundations-1718771, by Center for Brain-Inspired Computing, one of six centers in Joint University Microelectronics Program, a Semiconductor Research Corporation program sponsored by Defense Advanced Research Projects Agency, and by the EU Horizon 2020 research and innovation program under 646804-ERC-COG-BNYQ.

## REFERENCES

1. AHMED, A., AGHASI, A. & HAND, P. (2018) Blind deconvolutional phase retrieval via convex programming. *Advances in Neural Information Processing Systems*, pp. 10030–10040.
2. AHMED, A. & DEMANET, L. (2018) Leveraging diversity and sparsity in blind deconvolution. *IEEE Trans. Inf. Theory*, **64**, 3975–4000.
3. AHMED, A., RECHT, B. & ROMBERG, J. (2014) Blind deconvolution using convex programming. *IEEE Trans. Inf. Theory*, **60**, 1711–1732.
4. AHMED, A. & ROMBERG, J. (2015) Compressive multiplexing of correlated signals. *IEEE Trans. Inf. Theory*, **61**, 479–498.
5. ARIK, S. Ö. & KAHN, J. M. (2016) Direct-detection mode-division multiplexing in modal basis using phase retrieval. *Opt. Lett.*, **41**, 4265–4268.
6. BAHMANI, S. & ROMBERG, J. (2017) Phase retrieval meets statistical learning theory: a flexible convex relaxation. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, PMLR*, vol. 54. pp. 252–260.
7. BAHMANI, S. & ROMBERG, J. (2019) Solving equations of random convex functions via anchored regression. *Found. Comput. Math.*, **19**, 813–841.
8. BECKER, S. R., CANDÈS, E. J. & GRANT, M. C. (2011) Templates for convex cone problems with applications to sparse signal recovery. *Math. Program. Comput.*, **3**, 165.

9. BENDORY, T., EDIDIN, D. & ELДАР, Y. C. (2018) Blind phaseless short-time Fourier transform recovery. *IEEE Trans. Inf. Theory*, **66**, 3232–3241
10. BOUZIANE, R. & KILLEY, R. (2015) Blind symbol synchronization for direct detection optical OFDM using a reduced number of virtual subcarriers. *Opt. Express*, **23**, 6444–6454.
11. CAI, T. T., LI, X. & MA, Z. (2016) Optimal rates of convergence for noisy sparse phase retrieval via thresholded Wirtinger flow. *Ann. Stat.*, **44**, 2221–2251.
12. CANDÈS, E. & LI, X. (2014) Solving quadratic equations via PhaseLift when there are about as many equations as unknowns. *Found. Comput. Math.*, **14**, 1017–1026.
13. CANDÈS, E. & RECHT, B. (2009) Exact matrix completion via convex optimization. *Found. Comput. Math.*, **9**, 717–772.
14. CANDÈS, E., STROHMER, T. & VORONINSKI, V. (2013) PhaseLift: exact and stable signal recovery from magnitude measurements via convex programming. *Comm. Pure Appl. Math.*, **66**, 1241–1274.
15. CANDÈS, E. & TAO, T. (2010) The power of convex relaxation: near-optimal matrix completion. *IEEE Trans. Inf. Theory*, **56**, 2053–2080.
16. CANDÈS, E. J., LI, X. & SOLTANOLKOTABI, M. (2015) Phase retrieval from coded diffraction patterns. *Appl. Comp. Harm. Anal.*, **39**, 277–299.
17. CANDÈS, E. J., LI, X. & SOLTANOLKOTABI, M. (2015) Phase retrieval via Wirtinger flow: theory and algorithms. *IEEE Trans. Inf. Theory*, **61**, 1985–2007.
18. CHEN, Y. & CANDÈS, E. (2015) Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Advances in Neural Information Processing Systems 28*, (C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama & R. Garnett eds), pp. 739–747.
19. CHEN, Y., CHI, Y. & GOLDSMITH, A. J. (2015) Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Trans. Inf. Theory*, **61**, 4034–4059.
20. CHI, Y. (2016) Guaranteed blind sparse spikes deconvolution via lifting and convex optimization. *IEEE J. Sel. Topics Signal Process.*, **10**, 782–794.
21. COSSE, A. (2017) A note on the blind deconvolution of multiple sparse signals from unknown subspaces. *Wavelets and Sparsity XVII*, (Y. M. Lu, D. Van De Ville & M. Papadakis eds) vol. 10394, pp. 330–347. International Society for Optics and Photonics. doi: [10.1117/12.2272836](https://doi.org/10.1117/12.2272836).
22. DAVENPORT, M. A. & ROMBERG, J. (2016) An overview of low-rank matrix recovery from incomplete observations. *IEEE J. Sel. Topics Signal Process.*, **10**, 608–622.
23. DAVIS, C. & KAHAN, W. M. (1970) The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.*, **7**, 1–46.
24. DIRKSEN, S. (2016) Dimensionality reduction with subgaussian matrices: a unified theory. *Found. Comput. Math.*, **16**, 1367–1396.
25. ECKERT, R., TIAN, L. & WALLER, L. (2016) Algorithmic self-calibration of illumination angles in Fourier ptychographic microscopy. *Imaging and Applied Optics 2016*, pp. CT2D.3. Optical Society of America.
26. ELДАР, Y. C., SIDORENKO, P., MIXON, D. G., BAREL, S. & COHEN, O. (2015) Sparse phase retrieval from short-time Fourier measurements. *IEEE Signal Process. Lett.*, **22**, 638–642.
27. FOUCCART, S. & RAUHUT, H. (2013) *A Mathematical Introduction to Compressive Sensing*, vol. **1**. Basel: Birkhäuser.
28. GOLDSTEIN, T. & STUDER, C. (2017) Convex phase retrieval without lifting via PhaseMax. *Proceedings of the 34th International Conference on Machine Learning*. Vol. 70. pp. 1273–1281. Sydney, NSW, Australia: JMLR.org.
29. GOLUB, G. H. & VAN LOAN, C. F. (2012) *Matrix Computations*. Johns Hopkins Studies in the Mathematical Sciences. Baltimore, Maryland: Johns Hopkins University Press.
30. GROSS, D. (2011) Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inf. Theory*, **57**, 1548–1566.
31. HILLAR, C. J. & LIM, L.-H. (2013) Most tensor problems are NP-hard. *J. ACM*, **60**, 45.
32. HUANG, W. & HAND, P. (2018) Blind deconvolution by a steepest descent algorithm on a quotient manifold. *SIAM J. Imaging Sci.*, **11**, 2757–2785.

33. JAGANATHAN, K., ELДАР, Y. C. & HASSIBI, B. (2016) Phase retrieval: an overview of recent developments. *Optical Compressive Imaging*. ( B. RATON ed). FL: CRC Press, pp. 263–296.
34. JAGANATHAN, K., OYMAK, S. & HASSIBI, B. (2017) Sparse phase retrieval: uniqueness guarantees and recovery algorithms. *IEEE Trans. Signal Process.*, **65**, 2402–2410.
35. JUNGE, M. & ZENG, Q. (2013) Noncommutative Bennett and Rosenthal inequalities. *Ann. Probab.*, **41**, 4287–4316.
36. KESHAVAN, R. H., MONTANARI, A. & OH, S. (2010) Matrix completion from a few entries. *IEEE Trans. Inf. Theory*, **56**, 2980–2998.
37. KOLTCHINSKII, V., LOUNICI, K. & TSYBAKOV, A. B. (2011) Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Stat.*, **39**, 2302–2329.
38. LATAŁA, R. (2006) Estimates of moments and tails of Gaussian chaoses. *Ann. Probab.*, **34**, 2315–2331.
39. LECUE, G. & MENDELSON, S. (2015) Minimax rate of convergence and the performance of empirical risk minimization in phase recovery. *Electron. J. Probab.*, **20**, 1–29.
40. LEE, K., BAHMANI, S., ELДАР, Y. & ROMBERG, J. (2018) Phase retrieval of low-rank matrices. Presented at the 7th International Conference on Computational Harmonic Analysis.
41. LEE, K., KRAHMER, F. & ROMBERG, J. (2018) Spectral methods for passive imaging: nonasymptotic performance and robustness. *SIAM J. Imaging Sci.*, **11**, 2110–2164.
42. LEE, K., LI, Y., JUNGE, M. & BRESLER, Y. (2017) Blind recovery of sparse signals from subsampled convolution. *IEEE Trans. Inf. Theory*, **63**, 802–821.
43. LEE, K., TIAN, N. & ROMBERG, J. (2016) Fast and guaranteed blind multichannel deconvolution under a bilinear channel model. *Information Theory Workshop*. Cambridge, UK.
44. LI, X., LING, S., STROHMER, T. & WEI, K. (2018) Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Appl. Comput. Harmon. Anal.*
45. LI, X. & VORONINSKI, V. (2013) Sparse signal recovery from quadratic measurements via convex programming. *SIAM J. Math. Anal.*, **45**, 3019–3033.
46. LI, Y., LEE, K. & BRESLER, Y. (2018) Blind gain and phase calibration via sparse spectral methods. *IEEE Trans. Inf. Theory*, **65**, 3097–3123.
47. LING, S. & STROHMER, T. (2018) Self-calibration via linear least squares. *SIAM J. Imaging Sci.*, **11**, 252–292.
48. MANTZEL, W. & ROMBERG, J. (2015) Compressed subspace matching on the continuum. *Inf. Inference*, **4**, 79–107.
49. MANTZEL, W., ROMBERG, J. & SABRA, K. (2014) Round-robin multiple source localization. *J. Acoust. Soc. Am.*, **135**, 134–147.
50. MOULINES, E., DUHAMEL, P., CARDOSO, J.-F. & MAYRARGUE, S. (1995) Subspace methods for the blind identification of multichannel FIR filters. *IEEE Trans. Signal Process.*, **43**, 516–525.
51. NETRAPALLI, P., JAIN, P. & SANGHAVI, S. (2013) Phase retrieval using alternating minimization. *Advances in Neural Information Processing Systems 26*. (C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani & K. Q. Weinberger eds), pp. 2796–2804. Curran Associates, Inc. <http://papers.nips.cc/paper/5041-phase-retrieval-using-alternating-minimization.pdf>.
52. RECHT, B. (2011) A simpler approach to matrix completion. *J. Mach. Learn. Res.*, **12**, 3413–3430.
53. RECHT, B., FAZEL, M. & PARRILO, P. A. (2010) Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.*, **52**, 471–501.
54. SCHMIDL, T. M. & COX, D. C. (1997) Robust frequency and timing synchronization for OFDM. *IEEE Trans. Commun.*, **45**, 1613–1621.
55. SUN, J., QU, Q. & WRIGHT, J. (2018) A geometric analysis of phase retrieval. *Found. Comput. Math.*, **18**, 1131–1198.
56. TAN, Y. S. & VERSHYNIN, R. (2018) Phase retrieval via randomized Kaczmarz: theoretical guarantees. *Inf. Inference*, **8**, 97–123.
57. TIAN, N., BYUN, S.-H., SABRA, K. & ROMBERG, J. (2017) Multichannel myopic deconvolution in underwater acoustic channels via low-rank recovery. *J. Acoust. Soc. Am.*, **141**, 3337–3348.

58. TROPP, J. A. (2012) User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.*, **12**, 389–434.
59. VASWANI, N., NAYER, S. & ELDAR, Y. C. (2016) Low-rank phase retrieval. *IEEE Trans. Signal Process.*, **65**, 4059–4074.
60. VERSHYNIN, R. (2012) Introduction to the non-asymptotic analysis of random matrices. *Compressed Sensing: Theory and Applications*, ( Y. ELDAR, & G. KUTYNIOK eds), chap. 5. Cambridge, UK: Cambridge Univ. Press, pp. 210–268.
61. WALDSPURGER, I., D’ASPROMONT, A. & MALLAT, S. (2015) Phase recovery, MaxCut, and complex semidefinite programming. *Math. Program. Ser. A*, **149**, 47–81.
62. WANG, G., GIANNAKIS, G. B. & ELDAR, Y. C. (2017) Solving systems of random quadratic equations via truncated amplitude flow. *IEEE Trans. Inf. Theory*, **64**, 773–794.
63. WANG, L. & CHI, Y. (2016) Blind deconvolution from multiple sparse inputs. *IEEE Signal Process. Lett.*, **23**, 1384–1388.
64. WEDIN, P. (1972) Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, **12**, 99–111.
65. XU, G., LIU, H., TONG, L. & KAILATH, T. (1995) A least-squares approach to blind channel identification. *IEEE Trans. Signal Process.*, **43**, 2982–2993.
66. YANG, D., TANG, G. & WAKIN, M. B. (2016) Super-resolution of complex exponentials from modulations with unknown waveforms. *IEEE Trans. Inf. Theory*, **62**, 5809–5830.

### A. Expectations of symmetric Gaussian tensors

We repeatedly use the expectation of various tensor products of an i.i.d. Gaussian vector, which are summarized below. First, we consider the expectation of the fourth-order tensor product.

LEMMA A.1 Let  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Then

$$\mathbb{E} \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} = \sum_{j,k=1}^d (\mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k + \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_j \otimes \mathbf{e}_k + \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k \otimes \mathbf{e}_j),$$

where  $\mathbf{e}_j$  denotes the  $j$ th column of  $\mathbf{I}_d$  for  $j = 1, \dots, d$ .

PROOF OF LEMMA A.1 The expectation of  $\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}$  is written as

$$\begin{aligned} & \mathbb{E} \sum_{i_1, i_2, i_3, i_4=1}^d (\mathbf{e}_{i_1} \mathbf{e}_{i_1}^\top \otimes \mathbf{e}_{i_2} \mathbf{e}_{i_2}^\top \otimes \mathbf{e}_{i_3} \mathbf{e}_{i_3}^\top \otimes \mathbf{e}_{i_4} \mathbf{e}_{i_4}^\top) (\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}) \\ &= \sum_{i_1, i_2, i_3, i_4=1}^d \mathbb{E} g_{i_1} g_{i_2} g_{i_3} g_{i_4} (\mathbf{e}_{i_1} \otimes \mathbf{e}_{i_2} \otimes \mathbf{e}_{i_3} \otimes \mathbf{e}_{i_4}), \end{aligned}$$

where  $g_i$  denotes the  $i$ th entry of  $\mathbf{g}$  for  $i = 1, \dots, d$ . The proof completes by noting that all odd moments of a standard normal variable vanish.

The following lemma is a direct consequence of Lemma A.1.

LEMMA A.2 Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  and  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Then

$$\mathbb{E}(\mathbf{x}^\top \mathbf{g} \mathbf{g}^\top \mathbf{y}) \mathbf{g} \mathbf{g}^\top = (\mathbf{x}^\top \mathbf{y}) \mathbf{I}_d + \mathbf{x} \mathbf{y}^\top + \mathbf{y} \mathbf{x}^\top.$$

Next we consider the expectation of an 8-way tensor product applying to a fourth-order tensor product of a unit vector.

LEMMA A.3 Let  $\mathbf{x} \in \mathbb{S}^{d-1}$  and  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Then

$$\begin{aligned} \mathbb{E}(\mathbf{x}^\top \mathbf{g})^4 (\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}) &= 24(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}) \\ &+ 12 \sum_{l=1}^d (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l + \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{e}_l + \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \\ &\quad + \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l + \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{x} + \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x}) \\ &+ 3 \sum_{j,k=1}^d (\mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k + \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_j \otimes \mathbf{e}_k + \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k \otimes \mathbf{e}_j), \end{aligned} \quad (\text{A.1})$$

where  $\mathbf{e}_l$  denotes the  $l$ th column of  $\mathbf{I}_d$  for  $l = 1, \dots, d$ .

PROOF OF Lemma A.3 The expectation  $\mathbb{E}(\mathbf{x}^\top \mathbf{g})^4 (\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})$  is rewritten as

$$\begin{aligned} &\mathbb{E}(\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})(\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})^\top (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}) \\ &= \sum_{\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3, \mathbf{B}_4 \in \{\mathbf{P}_x, \mathbf{P}_{x^\perp}\}} \mathbb{E}(\mathbf{B}_1 \otimes \mathbf{B}_2 \otimes \mathbf{B}_3 \otimes \mathbf{B}_4)(\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})(\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})^\top (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}), \end{aligned} \quad (\text{A.2})$$

where  $\mathbf{P}_x$  and  $\mathbf{P}_{x^\perp}$  denote the orthogonal projection operators onto the subspace spanned by  $\mathbf{x}$  and its orthogonal complement, respectively.

If any of  $\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3, \mathbf{B}_4$  is different from the other three matrices, then the corresponding summand in (A.2) becomes zero since it has a factor that is an odd moment of  $\mathbf{x}^\top \mathbf{g} \sim \mathcal{N}(0, 1)$ . Therefore, it suffices to consider the following three cases.

**Case 1:**  $\mathbf{B}_1 = \mathbf{B}_2 = \mathbf{B}_3 = \mathbf{B}_4 = \mathbf{P}_x$ .

Since

$$\mathbf{P}_x \otimes \mathbf{P}_x \otimes \mathbf{P}_x \otimes \mathbf{P}_x = (\mathbf{x}\mathbf{x}^\top \otimes \mathbf{x}\mathbf{x}^\top \otimes \mathbf{x}\mathbf{x}^\top \otimes \mathbf{x}\mathbf{x}^\top) = (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x})(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x})^\top,$$

it follows that the corresponding summand is written as:

$$\begin{aligned} &(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}) \mathbb{E}(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x})^\top (\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})(\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g})^\top (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}) \\ &= \mathbb{E}(\mathbf{x}^\top \mathbf{g})^8 (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}) = 105(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}). \end{aligned} \quad (\text{A.3})$$

**Case 2:** Two of  $\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3, \mathbf{B}_4$  are  $\mathbf{P}_x$  and the other two matrices are  $\mathbf{P}_{x^\perp}$ .

First, we consider the sub-case where  $\mathbf{B}_1 = \mathbf{B}_2 = \mathbf{P}_x$  and  $\mathbf{B}_3 = \mathbf{B}_4 = \mathbf{P}_{x^\perp}$ . Since  $\mathbf{P}_{x^\perp} \mathbf{g}$  and  $\mathbf{x}^\top \mathbf{g}$  are independent, we can replace  $\mathbf{x}^\top \mathbf{g}$  by  $\mathbf{x}^\top \mathbf{g}'$  where  $\mathbf{g}'$  is an independent copy of  $\mathbf{g}$ . Then the corresponding summand is written as

$$\begin{aligned} &\mathbb{E}_{\mathbf{g}'}(\mathbf{g}'^\top \mathbf{x})^6 \mathbb{E}_{\mathbf{g}} \mathbf{P}_{x^\perp} \mathbf{g} \otimes \mathbf{P}_{x^\perp} \mathbf{g} \otimes \mathbf{x} \otimes \mathbf{x} = 15(\mathbf{P}_{x^\perp} \otimes \mathbf{P}_{x^\perp})(\mathbb{E}_{\mathbf{g}} \mathbf{g} \otimes \mathbf{g}) \otimes \mathbf{x} \otimes \mathbf{x} \\ &= 15(\mathbf{P}_{x^\perp} \otimes \mathbf{P}_{x^\perp}) \text{vec}(\mathbb{E}_{\mathbf{g}} \mathbf{g} \mathbf{g}^\top) \otimes \mathbf{x} \otimes \mathbf{x} = 15(\mathbf{P}_{x^\perp} \otimes \mathbf{P}_{x^\perp}) \text{vec}(\mathbf{I}_d) \otimes \mathbf{x} \otimes \mathbf{x} \\ &= 15 \text{vec}(\mathbf{P}_{x^\perp} \mathbf{I}_d \mathbf{P}_{x^\perp}) \otimes \mathbf{x} \otimes \mathbf{x} = 15 \text{vec}(\mathbf{I}_d - \mathbf{P}_x) \otimes \mathbf{x} \otimes \mathbf{x} \\ &= 15 \left( -\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} + \sum_{l=1}^d \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x} \right). \end{aligned}$$

The summands corresponding to the other sub-cases of Case 2 are calculated similarly, and the partial summation of (A.2) for Case 2 is written as

$$\begin{aligned}
 & -90\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} + 15 \sum_{l=1}^d (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l + \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{e}_l + \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \\
 & \quad + \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l + \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{x} + \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x}).
 \end{aligned} \tag{A.4}$$

**Case 3:**  $\mathbf{B}_1 = \mathbf{B}_2 = \mathbf{B}_3 = \mathbf{B}_4 = \mathbf{P}_{\mathbf{x}^\perp}$ .

Again by the independence between  $\mathbf{P}_{\mathbf{x}^\perp} \mathbf{g}$  and  $\mathbf{x}^\top \mathbf{g}$ , the corresponding summand is written as

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{g}'} (\mathbf{g}'^\top \mathbf{x})^4 \mathbb{E}_{\mathbf{g}} (\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) (\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}) \\
 & = 3(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) \mathbb{E}_{\mathbf{g}} (\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}).
 \end{aligned} \tag{A.5}$$

By plugging in the expression of  $\mathbb{E} \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}$  in Lemma A.1, the right-hand side of (A.5) is written as

$$\begin{aligned}
 & \underbrace{3(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) \sum_{j,k=1}^d \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k}_{(\S)} \\
 & \quad + \underbrace{3(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) \sum_{j,k=1}^d \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_j \otimes \mathbf{e}_k}_{(\S\S)} \\
 & \quad + \underbrace{3(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) \sum_{j,k=1}^d \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k \otimes \mathbf{e}_j}_{(\S\S\S)}.
 \end{aligned} \tag{A.6}$$

The first term (S) in (A.6) is rewritten as

$$\begin{aligned}
 (\S) & = (\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) [\text{vec}(\mathbf{I}_d) \otimes \text{vec}(\mathbf{I}_d)] \\
 & = [(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) \text{vec}(\mathbf{I}_d)] \otimes [(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp}) \text{vec}(\mathbf{I}_d)] \\
 & = \text{vec}(\mathbf{P}_{\mathbf{x}^\perp}) \otimes \text{vec}(\mathbf{P}_{\mathbf{x}^\perp}) = \text{vec}(\mathbf{I}_d - \mathbf{P}_{\mathbf{x}}) \otimes \text{vec}(\mathbf{I}_d - \mathbf{P}_{\mathbf{x}}) \\
 & = \text{vec}(\mathbf{P}_{\mathbf{x}}) \otimes \text{vec}(\mathbf{P}_{\mathbf{x}}) + \text{vec}(\mathbf{I}_d) \otimes \text{vec}(\mathbf{I}_d) - \text{vec}(\mathbf{I}_d) \otimes \text{vec}(\mathbf{P}_{\mathbf{x}}) - \text{vec}(\mathbf{P}_{\mathbf{x}}) \otimes \text{vec}(\mathbf{I}_d) \\
 & = \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} + \sum_{j,k=1}^d \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k - \sum_{l=1}^d (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l + \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x}).
 \end{aligned}$$

Similarly, (§§) and (§§§) are written as the sum of rank-1 tensors. Then applying these results to (A.6) provides

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{g}'}(\mathbf{g}'^\top \mathbf{x})^4 \mathbb{E}_{\mathbf{g}}(\mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp} \otimes \mathbf{P}_{\mathbf{x}^\perp})(\mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g} \otimes \mathbf{g}) \\
 &= 9 \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x} - 3 \sum_{l=1}^d (\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l + \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{e}_l + \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \\
 &\quad + \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{e}_l + \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{e}_l \otimes \mathbf{x} + \mathbf{e}_l \otimes \mathbf{e}_l \otimes \mathbf{x} \otimes \mathbf{x}) \\
 &+ 3 \sum_{j,k=1}^d (\mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k + \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_j \otimes \mathbf{e}_k + \mathbf{e}_j \otimes \mathbf{e}_k \otimes \mathbf{e}_k \otimes \mathbf{e}_j).
 \end{aligned} \tag{A.7}$$

The identity in (A.1) is then obtained by combining (A.3), (A.4) and (A.7) through (A.2).

**B. Moment and tail bounds of random matrices**

The following lemma, which provides a central moment bound on a standard normal variable, is a direct consequence of the Khintchine inequality (e.g. 60, Corollary 5.12).

LEMMA B.1 Let  $g \sim \mathcal{N}(0, 1)$ . Then there exists a numerical constant  $C$  such that

$$(\mathbb{E} |g|^p)^{1/p} \leq C \sqrt{p}, \quad \forall p \in \mathbb{N}.$$

We also use moment and tail bounds of random matrices in the spectral norm given by the noncommutative Rosenthal inequality [35, Theorem 0.4].

THEOREM B.2 (Noncommutative Rosenthal inequality [35, Theorem 0.4]). Let  $\mathbf{Y}_1, \dots, \mathbf{Y}_M$  be independent random matrices with zero-mean. Then there exists a numerical constant  $C$  such that

$$\left( \mathbb{E} \left\| \sum_{m=1}^M \mathbf{Y}_m \right\|^p \right)^{1/p} \leq C \left[ \sqrt{p} \left( \left\| \sum_{m=1}^M \mathbb{E} \mathbf{Y}_m \mathbf{Y}_m^* \right\|^{1/2} \vee \left\| \sum_{m=1}^M \mathbb{E} \mathbf{Y}_m^* \mathbf{Y}_m \right\|^{1/2} \right) \vee p \left( \sum_{m=1}^M \mathbb{E} \|\mathbf{Y}_m\|^p \right)^{1/p} \right]$$

for all  $1 \leq p < \infty$ .

Then the following lemma follows immediately from Theorem B.2.

LEMMA B.3 Let  $\mathbf{g}_1, \dots, \mathbf{g}_M \in \mathbb{R}^d$  be independent copies of  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ ,  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_M]^\top \in \mathbb{R}^M$ , and  $\nu \in (0, 1)$ . Then there exist numerical constants  $C_1, C_2 > 0$  such that

$$\left( \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M \lambda_m (\mathbf{g}_m \mathbf{g}_m^\top - \mathbf{I}_d) \right\|^p \right)^{1/p} \leq C_1 \|\boldsymbol{\lambda}\|_\infty \left[ M^{-1/2} \sqrt{pd} + M^{1/p-1} p(d+p) \right] \tag{B.1}$$

for all  $p \in \mathbb{N}$  and

$$\left\| \frac{1}{M} \sum_{m=1}^M \lambda_m (\mathbf{g}_m \mathbf{g}_m^\top - \mathbf{I}_d) \right\| \leq \delta$$

holds with probability  $1 - \nu$  provided

$$M \geq C_2 \left( \delta^{-1} \|\boldsymbol{\lambda}\|_\infty \vee \delta^{-2} \|\boldsymbol{\lambda}\|_\infty^2 \right) \left( d \log(M/\nu) \vee \log^2(M/\nu) \right). \tag{B.2}$$



**PROOF OF LEMMA B.3** We apply Theorem B.2 for  $Y_m = \lambda_m(\mathbf{g}_m\mathbf{g}_m^\top - \mathbf{I}_d)$  for  $m = 1, \dots, M$ . By the triangle inequality, we have

$$(\mathbb{E}\|Y_m\|^p)^{1/p} \leq \lambda_m + \lambda_m(\mathbb{E}\|\mathbf{g}_m\mathbf{g}_m^\top\|^p)^{1/p} = \lambda_m + \lambda_m(\mathbb{E}\|\mathbf{g}_m\|_2^{2p})^{1/p} \leq C_1\lambda_m(d+p).$$

Here the last step follows since:

$$\|\|\mathbf{g}\|_2 - \sqrt{d}\|_{L_{2p}} \leq C\sqrt{2p} \left\| \|\mathbf{g}\|_2 - \sqrt{d} \right\|_{\psi_2} \leq C'\sqrt{p},$$

where  $\|\cdot\|_{\psi_2}$  denotes the subgaussian norm. Therefore, we obtain

$$\left( \sum_{m=1}^M \mathbb{E}\|Y_m\|^p \right)^{1/p} \leq C_3\|\boldsymbol{\lambda}\|_\infty M^{1/p}(d+p). \tag{B.3}$$

Furthermore, the expectation of  $Y_m^2 = \lambda_m^2(\mathbf{g}_m\mathbf{g}_m^\top\mathbf{g}_m\mathbf{g}_m^\top - 2\mathbf{g}_m\mathbf{g}_m^\top + \mathbf{I}_d)$  is computed by using Lemma A.2 as  $\mathbb{E}Y_m^2 = \lambda_m^2(d+1)\mathbf{I}_d$ . Therefore, it follows that:

$$\left\| \sum_{m=1}^M \mathbb{E}Y_m^2 \right\|^{1/2} = \sqrt{d+1}\|\boldsymbol{\lambda}\|_2 \leq C_4\sqrt{Md}\|\boldsymbol{\lambda}\|_\infty. \tag{B.4}$$

Then (B.1) is obtained by plugging in (B.3) and (B.4) to Theorem B.2.

Next, by the Markov inequality, we have

$$\begin{aligned} \mathbb{P}\left( \left\| \frac{1}{M} \sum_{m=1}^M Y_m \right\| > \delta \right) &\leq \delta^{-p} \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M Y_m \right\|^p \\ &\leq C_1\delta^{-p} \|\boldsymbol{\lambda}\|_\infty^p \left[ M^{-1/2}\sqrt{pd} + M^{1/p-1}p(d+p) \right]^p. \end{aligned} \tag{B.5}$$

Let  $p = \log(M/\nu)$ . Then (B.2) implies that the right-hand side of (B.5) is upper-bounded by  $\nu$ . This completes the proof.

### C. Proof of Lemma 4.3

Let  $\phi := \angle(\mathbf{u}_0, \mathbf{u}_\#)$  and  $\psi := \angle(\mathbf{v}_0, \mathbf{v}_\#)$ . Then

$$\inf_{\theta \in [0, 2\pi)} \left\| \mathbf{u}_0\mathbf{v}_0^\top - e^{i\theta}\mathbf{u}_\#\mathbf{v}_\#^\top \right\|_F^2 = 2 - 2\cos\phi\cos\psi \leq 2 - 2\cos^2(\phi \vee \psi) = 2\sin^2(\phi \vee \psi).$$

Therefore, it suffices to show

$$\sin(\phi \vee \psi) = \sin\phi \vee \sin\psi \leq \frac{\delta}{\sqrt{2}}.$$

We will only show  $\sin\phi \leq \sqrt{\delta/2}$ . The derivation of the other part is essentially the same due to symmetry. Without loss of generality, we assume  $\|\mathbf{X}_\#\|_F = 1$  (or equivalently  $\sigma_\# = 1$ ).

Since  $\mathbf{X}_\#$  is a scalar multiple of the most dominant eigenvector of  $\mathbb{E}\mathbf{Y}$ , we use the Davis–Kahan theorem [23] to bound the error in estimating  $\mathbf{u}_\#$  as the dominant eigenvector of  $\mathbf{Y}$ . Among variations of the Davis–Kahan theorem, we use the version given in terms of the principal angle between two

subspaces. The following theorem states this result and is obtained by combining the argument of [29, Corollary 7.2.6] and the  $\sin \theta$  theorem for any unitarily invariant norm [64].

**THEOREM C.1** (Davis–Kahan  $\sin \theta$  theorem). Let  $\mathbf{A}, \mathbf{\Delta} \in \mathbb{C}^{n \times n}$  satisfy that  $\mathbf{A}$  and  $\mathbf{A} + \mathbf{\Delta}$  are positive semidefinite. Let  $\mathbf{Q} \in \mathbb{C}^{n \times r}$  (resp.  $\widehat{\mathbf{Q}} \in \mathbb{C}^{n \times r}$ ) denote the matrix whose columns are the eigenvectors of  $\mathbf{A}$  (resp.  $\mathbf{A} + \mathbf{\Delta}$ ) corresponding to the  $r$ -largest eigenvalues. Suppose that  $\lambda_r(\mathbf{A}) > \lambda_{r+1}(\mathbf{A})$ . If

$$\|\mathbf{\Delta}\| \leq \frac{\lambda_r(\mathbf{A}) - \lambda_{r+1}(\mathbf{A})}{5},$$

then

$$\sin \angle(\text{span}(\mathbf{Q}), \text{span}(\widehat{\mathbf{Q}})) \leq \frac{4\|\mathbf{\Delta}\|}{\lambda_r(\mathbf{A}) - \lambda_{r+1}(\mathbf{A})}.$$

To prove Lemma 4.3, we apply Theorem C.1 to  $\mathbf{A} = \mathbb{E} \mathbf{Y}$  and  $\mathbf{\Delta} = \mathbf{Y} - \mathbb{E} \mathbf{Y}$  with  $r = 1$ . Since

$$\mathbf{A} = \mathbf{u}_\# \mathbf{u}_\#^* + \left(1 + \frac{1}{M} \sum_{m=1}^M \xi_m\right) \mathbf{I}_{d_1},$$

it follows that:

$$\lambda_k(\mathbf{A}) = \lambda_k(\mathbf{u}_\# \mathbf{u}_\#^*) + 1 + \frac{1}{M} \sum_{m=1}^M \xi_m.$$

Therefore, we obtain

$$\lambda_1(\mathbf{A}) - \lambda_2(\mathbf{A}) = \lambda_1(\mathbf{u}_\# \mathbf{u}_\#^*) - \lambda_2(\mathbf{u}_\# \mathbf{u}_\#^*) = 1.$$

It remains to show

$$\|\mathbf{\Delta}\| \leq \frac{\delta}{4\sqrt{2}}. \quad (\text{C.1})$$

Let us first decompose  $\mathbf{\Delta}$  into its noise-free portion and the remainder as

$$\mathbf{\Delta} = \frac{1}{M} \sum_{m=1}^M |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \mathbf{a}_m \mathbf{a}_m^* - \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \mathbf{a}_m \mathbf{a}_m^* + \frac{1}{M} \sum_{m=1}^M \xi_m (\mathbf{a}_m \mathbf{a}_m^* - \mathbf{I}_{d_1}).$$

Then (C.1) is implied by

$$\left\| \frac{1}{M} \sum_{m=1}^M |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \mathbf{a}_m \mathbf{a}_m^* - \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \mathbf{a}_m \mathbf{a}_m^* \right\| \leq \frac{\delta}{8\sqrt{2}} \quad (\text{C.2})$$

and

$$\left\| \frac{1}{M} \sum_{m=1}^M \xi_m (\mathbf{a}_m \mathbf{a}_m^* - \mathbf{I}_{d_1}) \right\| \leq \frac{\delta}{8\sqrt{2}}. \quad (\text{C.3})$$

Indeed, by Lemma B.3, (24) implies that (C.3) holds with probability  $1 - \nu/2$  where  $\nu = M^{-\alpha}$ .

In the remainder of the proof, we show (22) implies (C.2) with probability  $1 - \nu/2$ . Let

$$\mathbf{Y}_m = \mathbf{Z}_m - \mathbb{E} \mathbf{Z}_m, \quad m = 1, \dots, M,$$

where

$$\mathbf{Z}_m = |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \mathbf{a}_m \mathbf{a}_m^*. \quad (\text{C.4})$$

Then (C.2) is written as

$$\left\| \frac{1}{M} \sum_{m=1}^M \mathbf{Y}_m \right\| \leq \frac{\delta}{8\sqrt{2}}. \quad (\text{C.5})$$

To show (C.5), we use the noncommutative Rosenthal inequality in Theorem B.2. By direct calculation, we obtain

$$\mathbb{E} \mathbf{Z}_m = \mathbf{u}_\# \mathbf{u}_\#^* + \mathbf{I}_{d_1}.$$

Next, by plugging in (C.4) into  $\mathbb{E} \mathbf{Z}_m^* \mathbf{Z}_m$ , we obtain

$$\mathbb{E} \mathbf{Z}_m^* \mathbf{Z}_m = \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{a}_m \mathbf{a}_m^* \mathbf{a}_m \mathbf{a}_m^*. \quad (\text{C.6})$$

By decomposing the right-hand side of (C.6) with  $\mathbf{P}_u + \mathbf{P}_{u^\perp} = \mathbf{I}_{d_1}$ ,  $\mathbb{E} \mathbf{Z}_m^* \mathbf{Z}_m$  is rewritten as

$$\mathbb{E} \mathbf{Z}_m^* \mathbf{Z}_m = \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \quad (\text{C.7a})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \quad (\text{C.7b})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \quad (\text{C.7c})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \quad (\text{C.7d})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \quad (\text{C.7e})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \quad (\text{C.7f})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#^\perp} \quad (\text{C.7g})$$

$$+ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{u_\#^\perp} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{u_\#}. \quad (\text{C.7h})$$

Since  $\mathbf{u}_\#^* \mathbf{a}_m$  and  $\mathbf{P}_{u_\#^\perp} \mathbf{a}_m$  are independent, which follows from  $\mathbf{a}_m \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{d_1})$ , we can substitute  $\mathbf{P}_{u_\#^\perp} \mathbf{a}_m$  by  $\mathbf{P}_{u_\#^\perp} \check{\mathbf{a}}_m$ , where  $\check{\mathbf{a}}_m$  is an independent copy of  $\mathbf{a}_m$ . For a standard complex Gaussian random variable  $\check{g} \sim \mathcal{CN}(0, 1)$ , we have

$$\mathbb{E} |\check{g}|^2 = 1, \quad \mathbb{E} |\check{g}|^4 = 2, \quad \mathbb{E} |\check{g}|^6 = 6, \quad \mathbb{E} |\check{g}|^8 = 24.$$

Therefore, by using these even-order moments of  $\mathcal{CN}(0, 1)$  together with the independence between  $\mathbf{a}_m$  and  $\check{\mathbf{a}}_m$ , we can compute (C.7a) to (C.7d) as follows:

$$\begin{aligned}\mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{\mathbf{u}_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{\mathbf{u}_\#} \mathbf{a}_m \mathbf{a}_m^* \mathbf{P}_{\mathbf{u}_\#} &= 48 \mathbf{P}_{\mathbf{u}_\#}, \\ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{\mathbf{u}_\#} \mathbf{a}_m \check{\mathbf{a}}_m^* \mathbf{P}_{\mathbf{u}_\#} \check{\mathbf{a}}_m \mathbf{a}_m^* \mathbf{P}_{\mathbf{u}_\#} &= 12(d_1 - 1) \mathbf{P}_{\mathbf{u}_\#}, \\ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{\mathbf{u}_\#} \check{\mathbf{a}}_m \mathbf{a}_m^* \mathbf{P}_{\mathbf{u}_\#} \mathbf{a}_m \check{\mathbf{a}}_m^* \mathbf{P}_{\mathbf{u}_\#} &= 12 \mathbf{P}_{\mathbf{u}_\#} \\ \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^4 |\mathbf{a}_m^* \mathbf{u}_\#|^4 \mathbf{P}_{\mathbf{u}_\#} \check{\mathbf{a}}_m \check{\mathbf{a}}_m^* \mathbf{P}_{\mathbf{u}_\#} \check{\mathbf{a}}_m \check{\mathbf{a}}_m^* \mathbf{P}_{\mathbf{u}_\#} &= 4(d_1 + 1) \mathbf{P}_{\mathbf{u}_\#}.\end{aligned}$$

Furthermore, each of the remaining summands (C.7e) to (C.7h) vanishes since it has a factor given as a central Gaussian moments of an odd order.

Applying the above results to (C.7) provides

$$\mathbb{E} \mathbf{Z}_m^* \mathbf{Z}_m = (12d_1 + 36) \mathbf{P}_{\mathbf{u}_\#} + (4d_1 + 16) \mathbf{P}_{\mathbf{u}_\#^\perp}.$$

Then, by the definition of  $\mathbf{Y}_m$ , we have

$$\mathbb{E} \mathbf{Y}_m^* \mathbf{Y}_m = \mathbb{E} \mathbf{Z}_m^* \mathbf{Z}_m - (\mathbb{E} \mathbf{Z}_m)^* (\mathbb{E} \mathbf{Z}_m) = (12d_1 + 32) \mathbf{P}_{\mathbf{u}_\#} + (4d_1 + 15) \mathbf{P}_{\mathbf{u}_\#^\perp}.$$

Therefore, for  $d_1 \geq 3$ , we have

$$\left\| \sum_{m=1}^M \mathbb{E} \mathbf{Y}_m \mathbf{Y}_m^* \right\|^{1/2} \vee \left\| \sum_{m=1}^M \mathbb{E} \mathbf{Y}_m^* \mathbf{Y}_m \right\|^{1/2} \leq C_1 \sqrt{Md_1}. \quad (\text{C.8})$$

Next we compute the  $p$ th moment of the spectral norm. The  $p$ th moment is considered as the norm in  $L_p$ . Then by the triangle inequality in  $L_p$ , we obtain

$$(\mathbb{E} \|\mathbf{Y}_m\|^p)^{1/p} \leq (\mathbb{E} \|\mathbf{Z}_m\|^p)^{1/p} + \|\mathbb{E} \mathbf{Z}_m\| \leq (\mathbb{E} \|\mathbf{Z}_m\|^p)^{1/p} + 2. \quad (\text{C.9})$$

Again by the triangle inequality, we obtain

$$\begin{aligned}(\mathbb{E} \|\mathbf{Z}_m\|^p)^{1/p} &= \left[ \mathbb{E} \left( |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \|\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}_m\|_2^2 + |\mathbf{b}_m^* \mathbf{v}_\#|^2 |\mathbf{a}_m^* \mathbf{u}_\#|^2 \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}_m\|_2^2 \right)^p \right]^{1/p} \\ &\leq \left( \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^{2p} |\mathbf{a}_m^* \mathbf{u}_\#|^{4p} \right)^{1/p} + \left( \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^{2p} |\mathbf{a}_m^* \mathbf{u}_\#|^2 \|\mathbf{P}_{\mathbf{u}_\#^\perp} \check{\mathbf{a}}_m\|_2^{2p} \right)^{1/p} \\ &\leq \left( \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^{2p} \right)^{1/p} \left( \mathbb{E} |\mathbf{a}_m^* \mathbf{u}_\#|^{4p} \right)^{1/p} + \left( \mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^{2p} \right)^{1/p} \left( \mathbb{E} |\mathbf{a}_m^* \mathbf{u}_\#|^{2p} \right)^{1/p} \left( \mathbb{E} \|\check{\mathbf{a}}_m\|_2^{2p} \right)^{1/p}.\end{aligned} \quad (\text{C.10})$$

Since  $\mathbf{a}_m^* \mathbf{u}_\# \sim \mathcal{CN}(0, 1)$  and  $\mathbf{b}_m^* \mathbf{v}_\# \sim \mathcal{CN}(0, 1)$ , by Lemma B.1, there exists a numerical constant  $C_2$  such that

$$(\mathbb{E} |\mathbf{a}_m^* \mathbf{u}_\#|^p)^{1/p} = (\mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^p)^{1/p} \leq C_2 \sqrt{p}.$$

Since  $2\|\check{\mathbf{a}}_m\|_2^2$  is a chi-square random variable of the degree-of-freedom  $2d_1$ , we obtain

$$\left( \mathbb{E} \|\check{\mathbf{a}}_m\|_2^{2p} \right)^{1/p} \leq C_3 (d_1 + p), \quad \forall p \geq 2.$$

Applying these upper estimates of the moments to (C.10) then to (C.9) provides

$$(\mathbb{E} \|\mathbf{Y}_m\|^p)^{1/p} \leq C_4(p^2 d_1 + p^3),$$

which implies

$$p \left( \sum_{m=1}^M \mathbb{E} \|\mathbf{Y}_m\|^p \right)^{1/p} \leq C_4 M^{1/p} (p^3 d_1 + p^4). \quad (\text{C.11})$$

By applying (C.8) and (C.11) to Theorem B.2, we obtain

$$\left( \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M \mathbf{Y}_m \right\|^p \right)^{1/p} \leq C_5 \left[ \sqrt{\frac{pd_1}{M}} + \frac{M^{1/p}(p^3 d_1 + p^4)}{M} \right] \quad (\text{C.12})$$

for all  $p \geq 2$  and  $d_1 \geq 3$ .

Finally, similar to [27, Proposition 7.11], we derive a tail bound from moment bounds. It follows from the Markov inequality that:

$$\mathbb{P} \left( \left\| \frac{1}{M} \sum_{m=1}^M \mathbf{Y}_m \right\| > \frac{\delta}{8\sqrt{2}} \right) \leq \left( \frac{8\sqrt{2}}{\delta} \right)^p \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M \mathbf{Y}_m \right\|^p. \quad (\text{C.13})$$

By plugging in (C.12) to (C.13), it follows that (C.5) holds with probability  $\nu$  provided that

$$C_6 \left[ \sqrt{\frac{pd_1}{M}} + \frac{M^{1/p}(p^3 d_1 + p^4)}{M} \right] \leq \delta \nu^{1/p}.$$

Then we set  $p = \log(M/\nu)$  so that (22) implies that (22) holds with probability  $1 - \nu/2$ . Therefore, the probability for violating (C.5) becomes  $\nu = M^{-\alpha}$ . This completes the proof.

#### D. Proof of Lemma 4.5

To simplify notation, let

$$\mathbf{Z}_m := \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle^2 \boldsymbol{\Phi}_m \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \boldsymbol{\Phi}_m^\top, \quad m = 1, \dots, M.$$

Then  $\boldsymbol{\Upsilon}$  is written as

$$\boldsymbol{\Upsilon} = \frac{1}{M} \sum_{m=1}^M (\mathbf{Z}_m + \underbrace{\xi_m \boldsymbol{\Phi}_m \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \boldsymbol{\Phi}_m^\top}_{(b)}). \quad (\text{D.1})$$

We derive the expectation of  $\boldsymbol{\Upsilon}$  in the following steps: first the expectation of the noise part (b) in (D.1) is computed as

$$\mathbb{E} \xi_m \boldsymbol{\Phi}_m \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \boldsymbol{\Phi}_m^\top = \xi_m \text{tr}(\widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top) \mathbf{I}_{d_1} = r \xi_m \mathbf{I}_{d_1}. \quad (\text{D.2})$$

Next we compute  $\mathbb{E} \mathbf{Z}_m$  by using Lemma A.1. Let  $\mathbf{x}_\sharp = \text{vec}(\mathbf{X}_\sharp)$  and  $\boldsymbol{\phi}_m = \text{vec}(\boldsymbol{\Phi}_m)$  for  $m = 1, \dots, M$ . Then  $\mathbf{Z}_m$  is rewritten as

$$\mathbf{Z}_m = (\text{tr} \otimes \mathbf{I}_{d_1}) \left[ (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \langle \boldsymbol{\phi}_m, \mathbf{x}_\sharp \rangle^2 \boldsymbol{\phi}_m \boldsymbol{\phi}_m^\top (\widehat{\mathbf{V}} \otimes \mathbf{I}_{d_1}) \right].$$

Since the partial trace operator is linear, the expectation of  $\mathbf{Z}_m$  is written as

$$\begin{aligned}\mathbb{E}\mathbf{Z}_m &= (\text{tr} \otimes \mathbf{I}_{d_1}) \left[ (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \mathbb{E} \langle \boldsymbol{\phi}_m, \mathbf{x}_\# \rangle^2 \boldsymbol{\phi}_m \boldsymbol{\phi}_m^\top (\widehat{\mathbf{V}} \otimes \mathbf{I}_{d_1}) \right] \\ &= (\text{tr} \otimes \mathbf{I}_{d_1}) \left[ (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) (2\mathbf{x}_\# \mathbf{x}_\#^\top + \|\mathbf{x}_\#\|_{\mathbb{F}}^2 \mathbf{I}_{d_1 d_2}) (\widehat{\mathbf{V}} \otimes \mathbf{I}_{d_1}) \right] \\ &= 2\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top + r \|\mathbf{x}_\#\|_{\mathbb{F}}^2 \mathbf{I}_{d_1},\end{aligned}\tag{D.3}$$

where the second identity follows from Lemma A.1. Then by combining (D.2) and (D.3), the expectation of  $\boldsymbol{\Upsilon}$  is written as

$$\mathbb{E}\boldsymbol{\Upsilon} = 2\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top + \left( r \|\mathbf{x}_\#\|_{\mathbb{F}}^2 + \frac{r}{M} \sum_{m=1}^M \xi_m \right) \mathbf{I}_{d_1}.\tag{D.4}$$

It follows from (25) that  $\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top$  in the right-hand side of (D.4) has rank- $r$  and its invariant space coincides with that of  $\mathbf{X}_\# \mathbf{X}_\#^\top = \mathbf{U}_\# \boldsymbol{\Sigma}_\#^2 \mathbf{U}_\#^\top$ . The inclusion of the former subspace to the latter is obvious from the construction. Furthermore, the rank of  $\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top$  is at most  $r$ . Indeed, the  $r$ th largest singular value of  $\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top$  satisfies

$$\begin{aligned}\sigma_r(\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top) &\geq \sigma_r(\mathbf{X}_\#)^2 \sigma_r(\widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{V}_\# \mathbf{V}_\#^\top) \\ &\geq \sigma_r(\mathbf{X}_\#)^2 \left( \sigma_r(\mathbf{V}_\# \mathbf{V}_\#^\top) - \|\mathbf{I}_{d_2} - \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top\| \right) \geq (1 - \delta_{\text{in}}) \sigma_r(\mathbf{X}_\#)^2,\end{aligned}$$

where the last step follows from (25). Therefore, we deduce that  $\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top$  and  $\mathbf{X}_\# \mathbf{X}_\#^\top$  have the same invariant subspace.

Recall that the columns of  $\mathbf{U}_0$  are the eigenvectors of  $\boldsymbol{\Upsilon}$  corresponding to the  $r$ -largest eigenvalues. Furthermore, the subspace spanned by the top  $r$  eigenvectors of  $\mathbb{E}\boldsymbol{\Upsilon}$  is the same to the column space of  $\mathbf{U}_\#$ . Therefore, the Davis–Kahan theorem (Theorem C.1) provides an upper bound for the estimation error measured by the principal angle between subspaces (the left-hand side of (28)). To this end, we apply Theorem C.1 to  $\mathbf{A} = \mathbb{E}\boldsymbol{\Upsilon}$  and  $\boldsymbol{\Delta} = \boldsymbol{\Upsilon} - \mathbb{E}\boldsymbol{\Upsilon}$  as shown below.

Since the spectral gap in  $\mathbf{A}$  satisfies

$$\lambda_r(\mathbf{A}) - \lambda_{r+1}(\mathbf{A}) = \lambda_r(\mathbb{E}\boldsymbol{\Upsilon}) - \lambda_{r+1}(\mathbb{E}\boldsymbol{\Upsilon}) = \lambda_r(2\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top) \geq 2(1 - \delta_{\text{in}}) [\sigma_r(\mathbf{X}_\#)]^2,$$

the error bound in (28) is obtained by Theorem C.1 provided that

$$\|\boldsymbol{\Delta}\| = \|\boldsymbol{\Upsilon} - \mathbb{E}\boldsymbol{\Upsilon}\| \leq \frac{(1 - \delta_{\text{in}}) \delta_{\text{out}} \sigma_r(\mathbf{X}_\#)^2}{2}.\tag{D.5}$$

By the triangle inequality, we obtain a sufficient condition for (D.5) given by

$$\left\| \frac{1}{M} \sum_{m=1}^M (\mathbf{Z}_m - \mathbb{E}\mathbf{Z}_m) \right\| \leq \frac{(1 - \delta_{\text{in}}) \delta_{\text{out}} \sigma_r(\mathbf{X}_\#)^2}{4}\tag{D.6}$$

and

$$\left\| \frac{1}{M} \sum_{m=1}^M \xi_m \left( \Phi_m \widehat{V} \widehat{V}^\top \Phi_m^\top - r \mathbf{I}_{d_1} \right) \right\| \leq \frac{(1 - \delta_{\text{in}}) \delta_{\text{out}} \sigma_r(\mathbf{X}_\#)^2}{4}. \tag{D.7}$$

In the remainder, we show that (D.6) and (D.6) hold with high probability when the conditions in (26) and (29) are satisfied. First, by Lemma B.3, it follows from (29) that (D.7) holds with probability  $1 - M^{-\alpha}/2$ . Then it remains to show that (D.6) holds with probability  $1 - M^{-\alpha}/2$  when (26) is satisfied. By the Markov inequality,

$$\begin{aligned} & \mathbb{P} \left( \left\| \frac{1}{M} \sum_{m=1}^M (\mathbf{Z}_m - \mathbb{E} \mathbf{Z}_m) \right\| > \frac{(1 - \delta_{\text{in}}) \delta_{\text{out}} \sigma_r(\mathbf{X}_\#)^2}{4} \right) \\ & \leq \left( \frac{4}{(1 - \delta_{\text{in}}) \delta_{\text{out}} \sigma_r(\mathbf{X}_\#)^2} \right)^p \cdot \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M (\mathbf{Z}_m - \mathbb{E} \mathbf{Z}_m) \right\|^p \end{aligned}$$

for any  $p > 0$ . Therefore, (D.6) holds with probability  $1 - M^{-\alpha}/2$  if

$$\underbrace{\left( \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M (\mathbf{Z}_m - \mathbb{E} \mathbf{Z}_m) \right\|^p \right)^{1/p}}_{(\ddagger)} \leq \frac{(1 - \delta_{\text{in}}) \delta_{\text{out}} \sigma_r(\mathbf{X}_\#)^2 M^{-\alpha/p}}{4}. \tag{D.8}$$

To get an upper estimate of  $(\ddagger)$  in (D.8), we apply the noncommutative Rosenthal inequality (Theorem B.2) to  $\mathbf{Y}_m = \mathbf{Z}_m - \mathbb{E} \mathbf{Z}_m$  for  $m = 1, \dots, M$ . The first step is to compute the expectation of  $\mathbf{Y}_m^2$  as follows: let  $\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3, \mathbf{Q}_4 \in \mathbb{R}^{d_1 \times d_2}$ . Note that each entry of  $\mathbf{Q}_1 \mathbf{Q}_2^\top \mathbf{Q}_3 \mathbf{Q}_4^\top$  is given as a linear combination of the entries of  $\text{vec}(\mathbf{Q}_1) \otimes \text{vec}(\mathbf{Q}_2) \otimes \text{vec}(\mathbf{Q}_3) \otimes \text{vec}(\mathbf{Q}_4)$ . Therefore, there exists a linear map  $\mathcal{R} : \mathbb{R}^{(d_1 d_2)^4} \rightarrow \mathbb{R}^{d_1 \times d_1}$  that satisfies

$$\mathcal{R}[\text{vec}(\mathbf{Q}_1) \otimes \text{vec}(\mathbf{Q}_2) \otimes \text{vec}(\mathbf{Q}_3) \otimes \text{vec}(\mathbf{Q}_4)] = \mathbf{Q}_1 \mathbf{Q}_2^\top \mathbf{Q}_3 \mathbf{Q}_4^\top.$$

We also define

$$\mathbf{T}_m := \langle \Phi_m, \mathbf{X}_\# \rangle^4 [\text{vec}(\Phi_m \widehat{V}) \otimes \text{vec}(\Phi_m \widehat{V}) \otimes \text{vec}(\Phi_m \widehat{V}) \otimes \text{vec}(\Phi_m \widehat{V})].$$

Then  $\mathbf{Z}_m^2$  is written as  $\mathbf{Z}_m^2 = \mathcal{R}(\mathbf{T}_m)$ . Since  $\text{vec}(\Phi_m \widehat{\mathbf{V}}) = (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \text{vec}(\Phi_m)$ , it follows that  $\mathbb{E}\mathbf{T}_m$  is written as:

$$\begin{aligned} \mathbb{E}\mathbf{T}_m &= [(\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \otimes (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \otimes (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \otimes (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1})] \mathbb{E}(\Phi_m^\top \text{vec}(\mathbf{X}_\#))^4 (\phi_m \otimes \phi_m \otimes \phi_m \otimes \phi_m) \\ &= 24 [\text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}})] \\ &\quad + 12 \sum_{k_1=1}^d \sum_{k_2=1}^d \left[ \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \right. \\ &\quad \quad + \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \\ &\quad \quad + \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \\ &\quad \quad + \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \\ &\quad \quad + \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \\ &\quad \quad \left. + \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{X}_\# \widehat{\mathbf{V}}) \right] \\ &\quad + 3 \sum_{j_1, k_1=1}^{d_1} \sum_{j_2, k_2=1}^{d_2} \left[ \text{vec}(\mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \right. \\ &\quad \quad + \text{vec}(\mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \\ &\quad \quad \left. + \text{vec}(\mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}}) \otimes \text{vec}(\mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}}) \right], \end{aligned}$$

where Lemma A.3 is used to compute  $\mathbb{E}\mathbf{T}_m$  in the second step. Also by the linearity of the map  $\mathcal{R}$ , it follows that:

$$\begin{aligned} \mathbb{E}\mathbf{Z}_m^2 &= \mathcal{R}(\mathbb{E}\mathbf{T}_m) \\ &= 24 \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \\ &\quad + 12 \|\mathbf{X}_\#\|_{\text{F}}^2 \sum_{l_1=1}^d \sum_{l_2=1}^d \left[ \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \mathbf{e}_{l_1} \widetilde{\mathbf{e}}_{l_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{l_2} \mathbf{e}_{l_1}^\top + \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{l_2} \mathbf{e}_{l_1}^\top \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{l_2} \mathbf{e}_{l_1}^\top \right. \\ &\quad \quad + \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{l_2} \mathbf{e}_{l_1}^\top \mathbf{e}_{l_1} \widetilde{\mathbf{e}}_{l_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top + \mathbf{e}_{l_1} \widetilde{\mathbf{e}}_{l_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{l_2} \mathbf{e}_{l_1}^\top \\ &\quad \quad \left. + \mathbf{e}_{l_1} \widetilde{\mathbf{e}}_{l_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \mathbf{e}_{l_1} \widetilde{\mathbf{e}}_{l_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top + \mathbf{e}_{l_1} \widetilde{\mathbf{e}}_{l_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{l_2} \mathbf{e}_{l_1}^\top \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \right] \\ &\quad + 3 \|\mathbf{X}_\#\|_{\text{F}}^4 \sum_{j_1, k_1=1}^{d_1} \sum_{j_2, k_2=1}^{d_2} \left[ \mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{j_2} \mathbf{e}_{j_1}^\top \mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{k_2} \mathbf{e}_{k_1}^\top \right. \\ &\quad \quad \left. + \mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{k_2} \mathbf{e}_{k_1}^\top \mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{k_2} \mathbf{e}_{k_1}^\top + \mathbf{e}_{j_1} \widetilde{\mathbf{e}}_{j_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{k_2} \mathbf{e}_{k_1}^\top \mathbf{e}_{k_1} \widetilde{\mathbf{e}}_{k_2}^\top \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \widetilde{\mathbf{e}}_{j_2} \mathbf{e}_{j_1}^\top \right]. \end{aligned}$$



After direct calculation, the above expression for  $\mathbb{E}\mathbf{Z}_m^2$  simplifies to

$$\begin{aligned} \mathbb{E}\mathbf{Z}_m^2 &= 24\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \\ &\quad + 12(2r + d_1 + 2) \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top + 12 \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \|\mathbf{X}_\# \widehat{\mathbf{V}}\|_{\mathbb{F}}^2 \mathbf{I}_{d_1} \\ &\quad + 3 \|\mathbf{X}_\#\|_{\mathbb{F}}^4 r(r + d_1 + 1) \mathbf{I}_{d_1}. \end{aligned} \tag{D.9}$$

Then, by combining (D.3) and (D.9), we obtain

$$\begin{aligned} \mathbb{E}\mathbf{Y}_m^2 &= \mathbb{E}\mathbf{Z}_m^2 - (\mathbb{E}\mathbf{Z}_m)^2 \\ &= 20\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top + 4(5r + 3d_1 + 6) \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top \\ &\quad + \left( 12 \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \|\mathbf{X}_\# \widehat{\mathbf{V}}\|_{\mathbb{F}}^2 + \|\mathbf{X}_\#\|_{\mathbb{F}}^4 r(2r + 3d_1 + 3) \right) \mathbf{I}_{d_1}. \end{aligned}$$

Therefore, the spectral norm of  $\mathbb{E}\mathbf{Y}_m^2$  is upper-bounded by

$$\|\mathbb{E}\mathbf{Y}_m^2\| \leq 20 \|\mathbf{X}_\#\|^4 + 4(5r + 3d_1 + 6) \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \|\mathbf{X}_\#\|^2 + 12 \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \|\mathbf{X}_\# \widehat{\mathbf{V}}\|_{\mathbb{F}}^2 + r(2r + 3d_1 + 3) \|\mathbf{X}_\#\|_{\mathbb{F}}^4.$$

Collecting the results for  $m = 1, \dots, M$  gives

$$\left\| \sum_{m=1}^M \mathbb{E}\mathbf{Y}_m^2 \right\|^{1/2} \leq Cr^{3/2} \sqrt{Md_1} \|\mathbf{X}_\#\|^2. \tag{D.10}$$

Moreover, by applying the triangle inequality in  $L_p$  twice to (D.3), we obtain

$$\begin{aligned} (\mathbb{E} \|\mathbf{Y}_m\|^p)^{1/p} &\leq \underbrace{\left[ \mathbb{E} \left( \langle \boldsymbol{\Phi}_m, \mathbf{X}_\# \rangle^2 \|\boldsymbol{\Phi}_m \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \boldsymbol{\Phi}_m^\top - r\mathbf{I}_{d_1}\| \right)^p \right]^{1/p}}_{(\text{†})} \\ &\quad + r \underbrace{\left[ \mathbb{E} \left( \langle \boldsymbol{\Phi}_m, \mathbf{X}_\# \rangle^2 - \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \right)^p \right]^{1/p}}_{(\text{††})} + \underbrace{2 \|\mathbf{X}_\# \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \mathbf{X}_\#^\top\|}_{(\text{†††})}. \end{aligned} \tag{D.11}$$

By the Cauchy–Schwarz inequality in  $L_2$ , the first term (†) on the right-hand side of (D.11) is upper-bounded by

$$(\text{†}) \leq \left( \mathbb{E} \langle \boldsymbol{\Phi}_m, \mathbf{X}_\# \rangle^{4p} \right)^{1/2p} \cdot \left( \mathbb{E} \|\boldsymbol{\Phi}_m \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \boldsymbol{\Phi}_m^\top - r\mathbf{I}_{d_1}\|^{2p} \right)^{1/2p}.$$

Since  $\langle \mathbf{X}_\#, \boldsymbol{\Phi}_m \rangle \sim \mathcal{N}(0, \|\mathbf{X}_\#\|_{\mathbb{F}}^2)$ , by Lemma B.1, we have

$$\left( \mathbb{E} \langle \boldsymbol{\Phi}_m, \mathbf{X}_\# \rangle^{4p} \right)^{1/2p} \leq Cp \|\mathbf{X}_\#\|_{\mathbb{F}}^2.$$

Then it follows from  $\text{vec}(\boldsymbol{\Phi}_m \widehat{\mathbf{V}}) = (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \boldsymbol{\phi}_m$  that

$$\mathbb{E} \text{vec}(\boldsymbol{\Phi}_m \widehat{\mathbf{V}}) \text{vec}(\boldsymbol{\Phi}_m \widehat{\mathbf{V}})^\top = \mathbb{E} (\widehat{\mathbf{V}}^\top \otimes \mathbf{I}_{d_1}) \boldsymbol{\phi}_m \boldsymbol{\phi}_m^\top (\widehat{\mathbf{V}} \otimes \mathbf{I}_{d_1}) \boldsymbol{\phi}_m = \widehat{\mathbf{V}}^\top \widehat{\mathbf{V}} \otimes \mathbf{I}_{d_1} = \mathbf{I}_{d_2 d_1},$$

which implies that  $\Phi_1 \widehat{\mathbf{V}}, \dots, \Phi_M \widehat{\mathbf{V}} \in \mathbb{R}^{d_1 \times r}$  are independent copies of a standard i.i.d. Gaussian matrix. Thus, Lemma B.3 implies

$$\left( \mathbb{E} \|\Phi_m \widehat{\mathbf{V}} \widehat{\mathbf{V}}^\top \Phi_m^\top - r \mathbf{I}_{d_1}\|^2 \right)^{1/2p} \leq C \left[ \sqrt{rpd_1} + r^{1/2p} p (d_1 + p) \right].$$

Then, by the triangle inequality in  $L_p$  and Lemma B.1, (44) is upper-bounded as

$$(44) \leq r \left( \mathbb{E} \langle \Phi_m, \mathbf{X}_\# \rangle^{2p} \right)^{1/p} + r \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \leq C' r p \|\mathbf{X}_\#\|_{\mathbb{F}}^2.$$

The last term is trivially upper-bounded by (44)  $\leq \|\mathbf{X}_\#\|_{\mathbb{F}}^2$ .

By collecting the above results, we obtain that the  $L_p$ -norm of  $\|\mathbf{Y}_m\|$  is upper-bounded by

$$\left( \mathbb{E} \|\mathbf{Y}_m\|^p \right)^{1/p} \leq C_1 \|\mathbf{X}_\#\|_{\mathbb{F}}^2 p \left( r + \sqrt{rpd_1} + r^{1/2p} p (d_1 + p) \right). \quad (\text{D.12})$$

Then, by applying (D.10) and (D.12) to Theorem B.2, we obtain that (4) in (D.8) is upper-bounded by

$$\begin{aligned} & \left( \mathbb{E} \left\| \frac{1}{M} \sum_{m=1}^M (\mathbf{Z}_m - \mathbb{E} \mathbf{Z}_m) \right\|^p \right)^{1/p} \\ & \leq C_3 \|\mathbf{X}_\#\|_{\mathbb{F}}^2 \left( r^{3/2} M^{-1/2} \sqrt{pd_1} + M^{1/p-1} p^2 r \left( r + \sqrt{rpd_1} + r^{1/2p} p (d_1 + p) \right) \right). \end{aligned}$$

Finally, we choose  $p = \log(M/M^{-\alpha}) = (\alpha + 1) \log M$ . Then (26) implies (D8). This completes the proof.

### E. Proof of Lemma 5.1

Since  $\widehat{\mathbf{X}}$  is a minimizer to (3), it satisfies

$$\text{Im} \langle \mathbf{X}_0, \widehat{\mathbf{X}} \rangle = 0. \quad (\text{E.1})$$

Then by (E.1) and (30) together with the fact that  $\text{rank}(\mathbf{X}_\#) = 1$ , we have

$$\text{Im} \langle \mathbf{X}_0, \mathbf{H} \rangle = 0,$$

where  $\mathbf{H} = \widehat{\mathbf{X}} - \mathbf{X}_\#$ . Let  $\mathcal{P}_{\mathbf{X}_\#}$  denote the orthogonal projection onto  $\mathbb{C}\mathbf{X}_\#$ , that is

$$\mathcal{P}_{\mathbf{X}_\#} : \mathbf{M} \mapsto \frac{\mathbf{X}_\# \langle \mathbf{X}_\#, \mathbf{M} \rangle}{\|\mathbf{X}_\#\|_{\mathbb{F}}^2}.$$

Then it follows that:

$$0 = |\text{Im} \langle \mathbf{X}_0, \mathbf{H} \rangle| \geq |\text{Im} \langle \mathcal{P}_{\mathbf{X}_\#}(\mathbf{X}_0), \mathbf{H} \rangle| - |\text{Im} \langle \mathbf{X}_0 - \mathcal{P}_{\mathbf{X}_\#}(\mathbf{X}_0), \mathbf{H} \rangle|,$$

which is rearranged as

$$|\text{Im} \langle \mathcal{P}_{\mathbf{X}_\#}(\mathbf{X}_0), \mathbf{H} \rangle| \leq |\text{Im} \langle \mathbf{X}_0 - \mathcal{P}_{\mathbf{X}_\#}(\mathbf{X}_0), \mathbf{H} \rangle|. \quad (\text{E..2})$$

By (30), the left-hand side of (E.2) is bounded from below as

$$\begin{aligned} |\operatorname{Im} \langle \mathcal{P}_{X_{\sharp}}(X_0), \mathbf{H} \rangle| &= \frac{|\operatorname{Im}(\langle X_0, X_{\sharp} \rangle \langle X_{\sharp}, \mathbf{H} \rangle)|}{\|X_{\sharp}\|_F^2} \\ &= \frac{\langle X_0, X_{\sharp} \rangle}{\|X_{\sharp}\|_F} \cdot \frac{|\operatorname{Im} \langle X_{\sharp}, \mathbf{H} \rangle|}{\|X_{\sharp}\|_F} \geq \sqrt{1 - \delta^2} \cdot \|X_0\|_F \cdot \frac{|\operatorname{Im} \langle X_{\sharp}, \mathbf{H} \rangle|}{\|X_{\sharp}\|_F}. \end{aligned}$$

Since the linear operator  $\iota : \mathbf{M} \mapsto \mathbf{M} - \mathcal{P}_{X_{\sharp}}(\mathbf{M})$  is self-adjoint and idempotent, the right-hand side of (E.2) is bounded from above as

$$\begin{aligned} |\operatorname{Im} \langle X_0 - \mathcal{P}_{X_{\sharp}}(X_0), \mathbf{H} \rangle| &= |\operatorname{Im} \langle X_0 - \mathcal{P}_{X_{\sharp}}(X_0), \mathbf{H} - \mathcal{P}_{X_{\sharp}}(\mathbf{H}) \rangle| \\ &\leq \|X_0 - \mathcal{P}_{X_{\sharp}}(X_0)\|_F \cdot \|\mathbf{H} - \mathcal{P}_{X_{\sharp}}(\mathbf{H})\|_F \leq \delta \|X_0\|_F \cdot \|\mathbf{H} - \mathcal{P}_{X_{\sharp}}(\mathbf{H})\|_F. \end{aligned}$$

Applying the above bounds to (E.2) completes the proof.

### F. Proof of Lemma 5.4

The following lemma provides a tail probability of the product of two jointly Gaussian variables.

LEMMA F.1 (A variation of [6, Lemma 5]). Let  $g_1, g_2$  be random variables that satisfy

$$\begin{bmatrix} g_1 \\ g_2 \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right).$$

Then for all  $t > 0$

$$\mathbb{P}(g_1 g_2 > t) \geq \frac{2}{\pi} \cos^{-1}\left(\frac{\sqrt{3 - \rho}}{2}\right) \exp\left(-\frac{2t}{1 + \rho}\right). \tag{F.1}$$

PROOF OF LEMMA F.1 Let  $w_1$  and  $w_2$  be independent copies of a standard normal random variable following  $\mathcal{N}(0, 1)$ . Then  $g_1$  and  $g_2$  are written as

$$g_1 = \sqrt{\frac{1 + \rho}{2}} w_1 + \sqrt{\frac{1 - \rho}{2}} w_2 \quad \text{and} \quad g_2 = \sqrt{\frac{1 - \rho}{2}} w_1 - \sqrt{\frac{1 + \rho}{2}} w_2.$$

With this representation, we have

$$\mathbb{P}(g_1 g_2 > t) = \mathbb{P}\left(\frac{1 + \rho}{2} w_1^2 - \frac{1 - \rho}{2} w_2^2 > t\right) = \mathbb{P}\left(\frac{\rho - 1}{2} + \frac{w_1^2}{w_1^2 + w_2^2} > \frac{t}{w_1^2 + w_2^2}\right).$$

Since  $w_1^2/(w_1^2 + w_2^2)$  and  $1/(w_1^2 + w_2^2)$  respectively depend only on the direction and the  $\ell_2$  norm of the standard normal random vector  $[w_1, w_2]^T$ , they are mutually independent. Furthermore,  $R = w_1^2 + w_2^2$  follows the exponential distribution with mean 1/2 and  $w_1/\sqrt{w_1^2 + w_2^2}$  is written as  $\cos \theta$  where  $\theta$  is a

uniform random variable on  $[0, 2\pi)$ . Then it follows that:

$$\begin{aligned} \mathbb{P}\left(\frac{\rho-1}{2} + \frac{w_1^2}{w_1^2 + w_2^2} > \frac{t}{w_1^2 + w_2^2}\right) &\geq \mathbb{P}\left(\frac{w_1^2}{w_1^2 + w_2^2} \geq \frac{3-\rho}{4} \text{ and } \frac{1+\rho}{4} > \frac{t}{w_1^2 + w_2^2}\right) \\ &= \mathbb{P}\left(\frac{w_1^2}{w_1^2 + w_2^2} \geq \frac{3-\rho}{4}\right) \mathbb{P}\left(w_1^2 + w_2^2 > \frac{4t}{1+\rho}\right) \\ &= \mathbb{P}\left(\cos^2 \theta \geq \frac{3-\rho}{4}\right) \mathbb{P}\left(R > \frac{4t}{1+\rho}\right). \end{aligned} \quad (\text{F.2})$$

The lower bound in (F.1) is obtained by computing the probabilities in (F.2).

We apply Lemma F.1 for  $g_1 = \langle \mathbf{X}_\#^\top, \Phi \rangle / \|\mathbf{X}_\#^\top\|_F$ ,  $g_2 = \langle \Phi, \mathbf{H} \rangle / \|\mathbf{H}\|_F$  and  $t = \tau'$ . Since the probability in (F.1) is a monotone increasing function in  $\rho$ , to get a lower bound on the tail probability, it suffices to compute a lower estimate of  $\rho$ .

Let  $\mathbf{X}_\# = \mathbf{U}_\# \Sigma_\# \mathbf{V}_\#^\top$  denote the SVD of  $\mathbf{X}_\#$ . Let  $\sigma_1, \dots, \sigma_r$  denote the singular values of  $\mathbf{X}_\#$  in the non-increasing order. Then  $\|\mathbf{X}_\#\|_* = \sum_{k=1}^r \sigma_k$ . By the triangle inequality, we have

$$\rho = \frac{\langle \mathbf{X}_\#, \mathbf{H} \rangle}{\|\mathbf{X}_\#^\top\|_F \|\mathbf{H}\|_F} \geq \underbrace{\frac{\langle \|\mathbf{X}_\#\|_* \mathbf{U}_\# \mathbf{V}_\#^\top, \mathbf{H} \rangle}{r \|\mathbf{X}_\#^\top\|_F \|\mathbf{H}\|_F}}_{(bb)} - \underbrace{\frac{|\langle r\mathbf{X}_\# - \|\mathbf{X}_\#\|_* \mathbf{U}_\# \mathbf{V}_\#^\top, \mathbf{H} \rangle|}{r \|\mathbf{X}_\#^\top\|_F \|\mathbf{H}\|_F}}_{(bbb)}. \quad (\text{F.3})$$

Note that, for all  $\mathbf{H} \in \mathcal{A}_\delta$ , the first summand (bb) is further bounded from below by

$$\frac{\langle \|\mathbf{X}_\#\|_* \mathbf{U}_\# \mathbf{V}_\#^\top, \mathbf{H} \rangle}{r \|\mathbf{X}_\#^\top\|_F \|\mathbf{H}\|_F} \geq -\frac{\|\mathbf{X}_\#\|_*}{r \|\mathbf{X}_\#^\top\|_F} \cdot \frac{\sqrt{r}\delta}{1-\lambda}.$$

The second term (bbb) can be upper-bounded by the Cauchy–Schwarz inequality with

$$\left\| \mathbf{X}_\# - \frac{\|\mathbf{X}_\#\|_* \mathbf{U}_\# \mathbf{V}_\#^\top}{r} \right\|_F \leq \frac{\sqrt{r}(\sigma_1 - \sigma_r)}{2}.$$

By plugging in the above estimates to (F.3), we obtain a sufficient condition for  $\rho \geq -0.9$  given by

$$\frac{\delta}{1-\lambda} \leq \frac{\sqrt{r} \|\mathbf{X}_\#^\top\|_F}{\|\mathbf{X}_\#\|_*} \cdot \left(0.9 - \frac{\sqrt{r}(\sigma_1(\mathbf{X}_\#) - \sigma_r(\mathbf{X}_\#))}{2 \|\mathbf{X}_\#^\top\|_F}\right). \quad (\text{F.4})$$

Here the right-hand side of (F.4) is no larger than  $(2.8 - \kappa)/2$ . Therefore, (14) implies that  $1 + \rho \geq 0.1$ . Then Lemma F.1 provides the lower bound in (40). This completes the proof.

### G. Proof of Lemma 5.5

Without loss of generality, we may assume  $\|\mathbf{X}_\#^\top\|_F = \|\mathbf{H}\|_F = 1$ . Since  $\mathcal{P}_T$  and  $\mathcal{P}_{T^\perp}$  are orthogonal projection operators onto corresponding subspaces, they are self-adjoint and idempotent linear operators. Therefore, it follows that:

$$\langle \Phi_m, \mathbf{H} \rangle = \langle \mathcal{P}_T(\Phi_m), \mathcal{P}_T(\mathbf{H}) \rangle + \langle \mathcal{P}_{T^\perp}(\Phi_m), \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle.$$

Then by Hölder's inequality, we obtain

$$\begin{aligned}
 \mathfrak{C}_M(\mathcal{A}_\delta) &= \mathbb{E} \sup_{\mathbf{H} \in \mathcal{A}_\delta} \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \langle \mathbf{X}_\sharp, \boldsymbol{\Phi}_m \rangle \langle \boldsymbol{\Phi}_m, \mathbf{H} \rangle \\
 &\leq \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle \right\|_{\mathbb{F}} \cdot \sup_{\mathbf{H} \in \mathcal{A}_\delta} \|\mathcal{P}_T(\mathbf{H})\|_{\mathbb{F}} \\
 &\quad + \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle \right\| \cdot \sup_{\mathbf{H} \in \mathcal{A}_\delta} \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_* \\
 &\leq \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle \right\|_{\mathbb{F}} \tag{G.1}
 \end{aligned}$$

$$+ \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle \right\| \cdot \frac{\sqrt{r}(1-\lambda+\delta)}{\lambda}, \tag{G.2}$$

where the last step follows from the expression of  $\mathcal{A}_\delta$  in (39).

The part in (G.1) is upper-bounded by

$$\begin{aligned}
 \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle \right\|_{\mathbb{F}} &\leq \sqrt{\mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle \right\|_{\mathbb{F}}^2} \\
 &= \sqrt{\frac{1}{M} \sum_{m=1}^M \mathbb{E} \|\mathcal{P}_T(\boldsymbol{\Phi}_m) \langle \boldsymbol{\Phi}_m, \mathbf{X}_\sharp \rangle\|_{\mathbb{F}}^2} = \sqrt{\mathbb{E} \|\mathcal{P}_T(\boldsymbol{\Phi}) \langle \boldsymbol{\Phi}, \mathbf{X}_\sharp \rangle\|_{\mathbb{F}}^2},
 \end{aligned}$$

where the first step follows from Jensen's inequality; the second step holds since  $(\epsilon_m)_{m=1}^M$  is a Rademacher sequence; the last step holds since  $\boldsymbol{\Phi}_1, \dots, \boldsymbol{\Phi}_M$  are independent copies of  $\boldsymbol{\Phi}$ . Indeed, since

$$\mathcal{P}_T(\boldsymbol{\Phi}) = U_\sharp U_\sharp^* \boldsymbol{\Phi} + (\mathbf{I}_{d_1} - U_\sharp U_\sharp^*) \boldsymbol{\Phi} V_\sharp V_\sharp^*,$$

it follows that:

$$\mathbb{E} \|\mathcal{P}_T(\boldsymbol{\Phi}) \langle \boldsymbol{\Phi}, \mathbf{X}_\sharp \rangle\|_{\mathbb{F}}^2 = \mathbb{E} \|U_\sharp U_\sharp^* \boldsymbol{\Phi}\|_{\mathbb{F}}^2 \langle \boldsymbol{\Phi}, \mathbf{X}_\sharp \rangle^2 + \|(\mathbf{I}_{d_1} - U_\sharp U_\sharp^*) \boldsymbol{\Phi} V_\sharp V_\sharp^*\|_{\mathbb{F}}^2 \langle \boldsymbol{\Phi}, \mathbf{X}_\sharp \rangle^2. \tag{G.3}$$

The first summand in the right-hand side of (G.3) is computed as

$$\begin{aligned}
 \mathbb{E} \|U_\sharp^T \boldsymbol{\Phi}\|_{\mathbb{F}}^2 \langle U_\sharp^T \boldsymbol{\Phi}, U_\sharp^T \mathbf{X}_\sharp \rangle^2 &= \text{tr} \left[ \mathbb{E} \langle \text{vec}(U_\sharp^T \boldsymbol{\Phi}), \text{vec}(U_\sharp^T \mathbf{X}_\sharp) \rangle^2 \text{vec}(U_\sharp^T \boldsymbol{\Phi}) \text{vec}(U_\sharp^T \boldsymbol{\Phi})^T \right] \\
 &= \text{tr} \left( 2 \text{vec}(U_\sharp^T \mathbf{X}_\sharp) \text{vec}(U_\sharp^T \mathbf{X}_\sharp)^T + \mathbf{I}_{rd_2} \right) = 2 + rd_2,
 \end{aligned}$$

where the second step follows from Lemma A.2 since  $\text{vec}(U_\sharp^T \boldsymbol{\Phi}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{rd_2})$ .

Let  $\Phi'$  be an independent copy of  $\Phi$ . Since  $(\mathbf{I}_{d_1} - U_{\#}U_{\#}^*)\Phi$  is independent of  $U_{\#}U_{\#}^*\Phi$ , the second summand in the right-hand side of (G.3) is written as

$$\begin{aligned} & \mathbb{E}\|(\mathbf{I}_{d_1} - U_{\#}U_{\#}^*)\Phi V_{\#}V_{\#}^*\|_{\mathbb{F}}^2 \langle U_{\#}U_{\#}^*\Phi', X_{\#} \rangle^2 \\ &= \mathbb{E}_{\Phi} \left\| \left( V_{\#}V_{\#}^* \otimes (\mathbf{I}_{d_1} - U_{\#}U_{\#}^*) \right) \text{vec}(\Phi) \right\|_2^2 \mathbb{E}_{\Phi'} \langle \Phi', X_{\#} \rangle^2 \\ &= \text{tr} \left( V_{\#}V_{\#}^* \otimes (\mathbf{I}_{d_1} - U_{\#}U_{\#}^*) \right) = r(d_1 - r). \end{aligned}$$

Therefore, we obtain

$$\mathbb{E}\|\mathcal{P}_T(\Phi) \langle \Phi, X_{\#} \rangle\|_{\mathbb{F}}^2 = r(d_1 + d_2 - r) + 2.$$

By Jensen's inequality, the expectation in (G.2) is upper-bounded by

$$\mathbb{E} \left\| \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^{\perp}}(\Phi_m) \langle \Phi_m, X_{\#} \rangle \right\| \leq \left( \mathbb{E} \left\| \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^{\perp}}(\Phi_m) \langle \Phi_m, X_{\#} \rangle \right\|^{2p} \right)^{1/2p}$$

for all  $p \in \mathbb{N}$ . Then we apply the noncommutative Rosenthal inequality (Theorem B.2) for

$$Y_m = \epsilon_m \mathcal{P}_{T^{\perp}}(\Phi_m) \langle \Phi_m, X_{\#} \rangle, \quad m = 1, \dots, M.$$

Since  $\mathcal{P}_T(X_{\#}) = X_{\#}$  and  $\mathcal{P}_T(\Phi_m)$  is independent from  $\mathcal{P}_{T^{\perp}}(\Phi_m)$ , it follows:

$$Y_m = \epsilon_m \mathcal{P}_{T^{\perp}}(\Phi_m) \langle \mathcal{P}_T(\Phi'_m), X_{\#} \rangle, \quad m = 1, \dots, M,$$

where  $\Phi'_1, \dots, \Phi'_M$  are independent copies of  $\Phi_1, \dots, \Phi_M$ . Furthermore, we have  $\mathbb{E} Y_m = \mathbf{0}$  for  $m = 1, \dots, M$ . By direct computation with Lemma A.2, we obtain

$$\mathbb{E} Y_m Y_m^{\top} = \text{tr}(\mathbf{P}_{V_{\#}^{\perp}}) \mathbf{P}_{U_{\#}^{\perp}} \quad \text{and} \quad \mathbb{E} Y_m^{\top} Y_m = \text{tr}(\mathbf{P}_{U_{\#}^{\perp}}) \mathbf{P}_{V_{\#}^{\perp}}, \quad m = 1, \dots, M.$$

Therefore, we obtain

$$\left\| \sum_{m=1}^M \mathbb{E} Y_m Y_m^{\top} \right\|^{1/2} \vee \left\| \sum_{m=1}^M \mathbb{E} Y_m^{\top} Y_m \right\|^{1/2} \leq \sqrt{M(d_1 + d_2)}.$$

Next we derive an upper bound for  $\left( \sum_{m=1}^M \mathbb{E} \|Y_m\|^{2p} \right)^{1/2p}$ , which coincides with  $M^{1/2p} (\mathbb{E} \|Y_m\|^{2p})^{1/2p}$  for any  $m \in \{1, \dots, M\}$ . Since  $\mathcal{P}_T(\Phi_m)$  and  $\mathcal{P}_{T^{\perp}}(\Phi_m)$  are independent, it follows that the spectral norm of  $Y_m$  satisfies:

$$\mathbb{E} \|Y_m\|^{2p} \leq \left( \mathbb{E} \|\mathcal{P}_{T^{\perp}}(\Phi_m)\|^{2p} \right) \cdot \left( \mathbb{E} |\langle \mathcal{P}_T(\Phi_m), X_{\#} \rangle|^{2p} \right) \leq (C\sqrt{p})^{2p} \mathbb{E} \|\Phi_m\|^{2p},$$

where the last inequality follows from the fact that  $\langle \mathcal{P}_T(\Phi_m), X_{\#} \rangle \sim \mathcal{N}(0, 1)$  satisfies:

$$\mathbb{E} |\langle \mathcal{P}_T(\Phi_m), X_{\#} \rangle|^{2p} \leq (C\sqrt{p})^{2p}.$$

It remains to get an upper bound on  $\mathbb{E} \|\Phi_m\|^{2p}$ . Note that  $\|\Phi_m\|^2 = \|\Phi_m^{\top} \Phi_m\|$  where  $\Phi_m^{\top} \Phi_m$  follows the Wishart distribution. Without loss of generality, we may assume  $d_1 \leq d_2$  (otherwise we consider  $\Phi_m \Phi_m^{\top}$

instead of  $\Phi_m^\top \Phi_m$ ). Then Lemma B.3 implies

$$\left(\mathbb{E}\|\Phi_m^\top \Phi_m\|^p\right)^{1/p} \leq d_1 + C_1 \left(\sqrt{pd_1 d_2} + pd_1^{1/p} (d_2 + p)\right). \tag{G.4}$$

Indeed, (G.4) is obtained by Lemma B.3 and the triangle inequality in the Banach space of random variables  $L_p(\Omega, \mu)$ . Note that  $\Phi_m^\top \Phi_m$  is written as

$$\Phi_m^\top \Phi_m = \sum_{k=1}^{d_1} \mathbf{g}_k \mathbf{g}_k^\top,$$

where  $\mathbf{g}_1, \dots, \mathbf{g}_{d_1}$  are independent copies of  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_2})$ . Then it follows by Lemma B.3 that:

$$\left(\mathbb{E}\left\|\frac{1}{d_1} \Phi_m^\top \Phi_m - \mathbf{I}_{d_2}\right\|^p\right)^{1/p} \leq C_1 \left(d_1^{-1/2} \sqrt{pd_2} + d_1^{1/p-1} p (d_2 + p)\right),$$

which, together with the triangle inequality and the homogeneity of  $L_p$ -norm, implies (G.4). Then taking the square root on both sides of (G.4) gives

$$\left(\mathbb{E}\|\Phi_m\|^{2p}\right)^{1/2p} \leq C_2 \left(\sqrt{d_1} + (pd_1 d_2)^{1/4} + \sqrt{pd_1^{1/2p}} \sqrt{d_2} + pd_1^{1/2p}\right).$$

By collecting the above estimates, Theorem B.2 implies

$$\begin{aligned} & \left(\mathbb{E}\left\|\frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\Phi_m) \langle \Phi_m, \mathbf{X}_\# \rangle\right\|^{2p}\right)^{1/2p} \\ & \leq C_3 \left(\sqrt{p(d_1 + d_2)} + M^{1/2p-1/2} p^{3/2} \left(\sqrt{d_1} + (pd_1 d_2)^{1/4} + \sqrt{pd_1^{1/2p}} \sqrt{d_2} + pd_1^{1/2p}\right)\right). \end{aligned}$$

For the brevity, let  $d = d_1 + d_2$ . Let us choose  $p = 1 \vee \log d$ . Then  $1/2p - 1/2 \leq 0$ . Since  $M \geq d$ , we obtain

$$M^{1/2p-1/2} \leq d^{1/2p-1/2} \leq \frac{d^{1/2 \log d}}{\sqrt{d}} \leq C_4 d^{-1/2}.$$

Furthermore, we have

$$\sqrt{d_1} + (pd_1 d_2)^{1/4} + \sqrt{pd_1^{1/2p}} \sqrt{d_2} + p^{3/2} d_1^{1/2p} d_2^{1/2p} \leq C_5 \left(\sqrt{d \log d} + \log^{3/2} d\right).$$

Combining the above estimates provides

$$\begin{aligned} \mathbb{E}\left\|\frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\Phi_m) \langle \Phi_m, \mathbf{X}_\# \rangle\right\| & \leq \left(\mathbb{E}\left\|\frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\Phi_m) \langle \Phi_m, \mathbf{X}_\# \rangle\right\|^{2p}\right)^{1/2p} \\ & \leq C_6 \sqrt{(d_1 + d_2) \log(d_1 + d_2)}. \end{aligned}$$

Then the upper bound in (41) is obtained by applying the above estimates to (G.1) and (G.2).

## H. Proof of Lemma 5.6

The event is determined by a 1-homogeneous equation in  $\mathbf{H}$  and  $\mathbf{X}_\#$ . Therefore, without loss of generality, we may assume that  $\|\mathbf{H}\|_F = \|\mathbf{X}_\#\|_F = 1$ . Then  $\mathbf{X}_\#$  is written as  $\mathbf{u}_\# \mathbf{v}_\#^*$  with  $\|\mathbf{u}_\#\|_2 = \|\mathbf{v}_\#\|_2 = 1$ .

First, we decompose  $\text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{H} \mathbf{b})$  as

$$\begin{aligned} \text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{H} \mathbf{b}) &= \text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#} \mathbf{b}) + \text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#} \mathbf{b}) \\ &\quad + \text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}) + \text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}). \end{aligned}$$

By plugging in  $\mathbf{P}_{\mathbf{u}_\#} = \mathbf{u}_\# \mathbf{u}_\#^*$  and  $\mathbf{P}_{\mathbf{v}_\#} = \mathbf{v}_\# \mathbf{v}_\#^*$  to the above identity, we rewrite  $\text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{H} \mathbf{b})$  as

$$\begin{aligned} \text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{H} \mathbf{b}) &= |\mathbf{v}_\#^* \mathbf{b}|^2 |\mathbf{u}_\#^* \mathbf{a}|^2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + |\mathbf{v}_\#^* \mathbf{b}|^2 |\mathbf{u}_\#^* \mathbf{a}| \text{Re} \left( \frac{\mathbf{u}_\#^* \mathbf{a}}{|\mathbf{u}_\#^* \mathbf{a}|} \cdot \mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\# \right) \\ &\quad + |\mathbf{u}_\#^* \mathbf{a}|^2 |\mathbf{v}_\#^* \mathbf{b}| \text{Re} \left( \frac{\overline{\mathbf{v}_\#^* \mathbf{b}}}{|\mathbf{v}_\#^* \mathbf{b}|} \cdot \mathbf{u}_\#^* \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b} \right) + |\mathbf{u}_\#^* \mathbf{a}| |\mathbf{v}_\#^* \mathbf{b}| \text{Re} \left( \frac{\mathbf{u}_\#^* \mathbf{a}}{|\mathbf{u}_\#^* \mathbf{a}|} \cdot \frac{\overline{\mathbf{v}_\#^* \mathbf{b}}}{|\mathbf{v}_\#^* \mathbf{b}|} \cdot \mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b} \right). \end{aligned}$$

The following facts follow from the assumption that  $\mathbf{a} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{d_1})$  and  $\mathbf{b} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{d_2})$  are mutually independent:

1.  $|\mathbf{u}_\#^* \mathbf{a}|$ ,  $|\mathbf{v}_\#^* \mathbf{b}|$ ,  $\overline{\mathbf{u}_\#^* \mathbf{a}}/|\mathbf{u}_\#^* \mathbf{a}|$ ,  $\overline{\mathbf{v}_\#^* \mathbf{b}}/|\mathbf{v}_\#^* \mathbf{b}|$ ,  $\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}$ , and  $\mathbf{P}_{\mathbf{v}_\#} \mathbf{b}$  are independent random variables.
2.  $|\mathbf{u}_\#^* \mathbf{a}|$  and  $|\mathbf{v}_\#^* \mathbf{b}|$  follow the Rayleigh distribution with scale parameter 1.
3.  $\overline{\mathbf{u}_\#^* \mathbf{a}}/|\mathbf{u}_\#^* \mathbf{a}|$  and  $\overline{\mathbf{v}_\#^* \mathbf{b}}/|\mathbf{v}_\#^* \mathbf{b}|$  follow the uniform distribution on the set of complex number of the unit modulus.

Furthermore, due to the rotation invariance of the Gaussian distribution,  $(\overline{\mathbf{u}_\#^* \mathbf{a}}/|\mathbf{u}_\#^* \mathbf{a}|) \mathbf{P}_{\mathbf{u}_\#} \mathbf{a}$  has the same distribution with  $\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}$ . Similarly,  $(\overline{\mathbf{v}_\#^* \mathbf{b}}/|\mathbf{v}_\#^* \mathbf{b}|) \mathbf{P}_{\mathbf{v}_\#} \mathbf{b}$  and  $\mathbf{P}_{\mathbf{v}_\#} \mathbf{b}$  have the same distribution.

Combining the above facts, we obtain that  $\text{Re}(\mathbf{b}^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a} \mathbf{a}^* \mathbf{H} \mathbf{b})$  has the same distribution with

$$x := r_1^2 r_2^2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_1 r_2^2 \text{Re}(\mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#) + r_1^2 r_2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}) + r_1 r_2 \text{Re}(\mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}),$$

where  $r_1, r_2, \mathbf{a}, \mathbf{b}$  are independent and  $r_1, r_2 \sim \text{Rayleigh}(1)$ .

Now it suffices to compute the probability of the event  $\mathcal{E}$  defined by

$$\mathcal{E} := \{x \geq \tau'\}.$$

For positive constants  $\alpha, \beta$ , we define another event  $\mathcal{E}_0$  by

$$\mathcal{E}_0 := \{\alpha \leq r_1 \leq \beta, \alpha \leq r_2 \leq \beta\}.$$

For example, we may set  $\alpha = 0.9$  and  $\beta = 1.1$ . Then  $\mathbb{P}(\mathcal{E}_0) \geq 0.12$ .

Let  $z_1, z_2, z_3$  be random variables defined by

$$z_1 := \text{Re}(\mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#), \quad z_2 := \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}) \quad \text{and} \quad z_3 := \text{Re}(\mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}).$$



Since  $\mathbf{a} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{d_1})$ , if  $\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\# \neq \mathbf{0}$ , then it follows that  $\mathbf{a}^* \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\# \sim \mathcal{CN}(0, \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#\|_{\mathbb{F}}^2)$  and its real part  $z_1$  follows  $\mathcal{N}(0, \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#\|_{\mathbb{F}}^2/2)$ . Otherwise,  $\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\# = \mathbf{0}$  implies  $z_1 = 0$ . Similarly,  $z_2 \sim \mathcal{N}(0, \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{H}^* \mathbf{P}_{\mathbf{u}_\#}\|_{\mathbb{F}}^2/2)$  if  $\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{H}^* \mathbf{u}_\# \neq \mathbf{0}$ ;  $z_2 = 0$  otherwise. By the independence between  $\mathbf{a}$  and  $\mathbf{b}$ , it follows that  $z_1 + z_2 \sim \mathcal{N}(0, \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#\|_{\mathbb{F}}^2 + \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{H}^* \mathbf{u}_\#\|_{\mathbb{F}}^2/2)$  if  $\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#\|_{\mathbb{F}}^2 + \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{H}^* \mathbf{u}_\#\|_{\mathbb{F}}^2 > 0$ ;  $z_1 + z_2 = 0$  otherwise. In both cases,  $z_1 + z_2$  has a symmetric distribution, that is  $z_1 + z_2$  is equivalent to  $-(z_1 + z_2)$  in distribution.

Furthermore, we can rewrite  $z_3$  as a Gaussian bilinear form, i.e.

$$z_3 = \tilde{\mathbf{a}}^\top \mathbf{Q} \tilde{\mathbf{b}}$$

for

$$\mathbf{Q} = \frac{1}{2} \begin{bmatrix} \text{Re}(\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#^\perp) & -\text{Im}(\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#^\perp) \\ \text{Im}(\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#^\perp) & \text{Re}(\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#^\perp) \end{bmatrix},$$

where

$$\tilde{\mathbf{a}} = \sqrt{2} \begin{bmatrix} \text{Re}(\mathbf{a}) \\ \text{Im}(\mathbf{a}) \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2d_1}), \quad \tilde{\mathbf{b}} = \sqrt{2} \begin{bmatrix} \text{Re}(\mathbf{b}) \\ \text{Im}(\mathbf{b}) \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2d_2}).$$

It follows that  $z_3$  has a symmetric distribution. Furthermore, since  $z_3$  is a Gaussian bilinear form, it has a mixed subexponential-subgaussian tail given by

$$\mathbb{P}(|z_3| \geq t) \leq C \exp \left[ -\frac{1}{C} \left( \frac{t^2}{\|\mathbf{Q}\|_{\mathbb{F}}^2} \wedge \frac{t}{\|\mathbf{Q}\|} \right) \right], \quad \forall t > 0 \tag{H.1}$$

for a numerical constant  $C$ . Latała [38] showed that this tail bound is tight with an analogous lower bound given by

$$\mathbb{P}(|z_3| \geq t) \geq \frac{1}{C} \exp \left[ -C \left( \frac{t^2}{\|\mathbf{Q}\|_{\mathbb{F}}^2} \wedge \frac{t}{\|\mathbf{Q}\|} \right) \right], \quad \forall t > 0.$$

By direct calculation, we obtain

$$\|\mathbf{Q}\|_{\mathbb{F}} = \frac{1}{\sqrt{2}} \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#^\perp\|_{\mathbb{F}}$$

and

$$\|\mathbf{Q}\| = \frac{1}{2} \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{v}_\#^\perp\|.$$

Now we are ready to derive a lower bound on the probability of the event  $\mathcal{E}$  using the aforementioned properties  $z_1, z_2, z_3$ . It follows from the definition of the conditional probability that

$$\frac{\mathbb{P}(\mathcal{E})}{\mathbb{P}(\mathcal{E}_0)} \geq \frac{\mathbb{P}(\mathcal{E} \cap \mathcal{E}_0)}{\mathbb{P}(\mathcal{E}_0)} = \mathbb{P}(\mathcal{E} | \mathcal{E}_0) = \mathbb{P} \left( r_1 r_2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 + z_3 \geq \frac{\tau'}{r_1 r_2} \mid \mathcal{E}_0 \right). \tag{H.2}$$

As we choose  $\alpha < \beta$  as numerical constants,  $\mathbb{P}(\mathcal{E}_0)$  is another numerical constant. It remains to show that the lower bound in (H.2) is larger than a numerical constant. We consider the two complementary scenarios below.

**Case 1:** First, we consider the case when  $\mathbf{H}$  satisfies

$$\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp}\| \leq \zeta \quad (\text{H.3})$$

for some constant  $0 < \zeta < 1$ , which we will specify later.

Let  $\tau'' > 0$ . Then by the inclusion–exclusion principle, the right-hand side of (H.2) is lower-bounded by

$$\begin{aligned} & \mathbb{P}\left(r_1 r_2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 + z_3 \geq \frac{\tau'}{r_1 r_2} \mid \mathcal{E}_0\right) \\ & \geq \mathbb{P}\left(r_1 r_2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{r_1 r_2} \mid \mathcal{E}_0\right) - \mathbb{P}\left(z_3 < -\frac{\tau''}{r_1 r_2} \mid \mathcal{E}_0\right). \end{aligned} \quad (\text{H.4})$$

In the sequel, we will use the fact that for then the tail probability is

$$g \sim \mathcal{N}(0, \zeta^2) \implies \mathbb{P}(g > t) \text{ is increasing in } \zeta \text{ and decreasing in } t \geq 0. \quad (\text{H.5})$$

The following cases on the sign of  $\operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#)$  have to be distinguished. First, suppose that  $\operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) \leq 0$ . Conditioned on  $r_1$  and  $r_2$ , the random variable  $r_2 z_1 + r_1 z_2$  becomes a Gaussian and invoking (H.5) yields

$$\begin{aligned} & \mathbb{P}\left(r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{r_1 r_2} - r_1 r_2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) \mid \mathcal{E}_0, r_1, r_2\right) \\ & \geq \mathbb{P}\left(\alpha(z_1 + z_2) \geq \frac{\tau' + \tau''}{\alpha^2} - \beta^2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) \mid \mathcal{E}_0, r_1, r_2\right). \end{aligned}$$

Since the right-hand side of the above inequality is independent of  $r_1$  and  $r_2$ , we can conclude that

$$\mathbb{P}\left(r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{r_1 r_2} - r_1 r_2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) \mid \mathcal{E}_0\right) \geq \mathbb{P}\left(\alpha(z_1 + z_2) \geq \frac{\tau' + \tau''}{\alpha^2} - \beta^2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#)\right).$$

Furthermore,  $z_3$  is symmetric and we obtain an upper estimate of the tail probability of  $z_3$  in (H.4) given by

$$\mathbb{P}\left(z_3 < -\frac{\tau''}{r_1 r_2} \mid \mathcal{E}_0\right) \leq \mathbb{P}\left(z_3 < -\frac{\tau''}{\beta^2}\right) = \frac{1}{2} \mathbb{P}\left(|z_3| \geq \frac{\tau''}{\beta^2}\right).$$

By combining the above bounds, the lower estimate in (H.4) is further bounded from below by

$$\begin{aligned} & \mathbb{P}\left(r_1 r_2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{r_1 r_2} \mid \mathcal{E}_0\right) - \mathbb{P}\left(z_3 < -\frac{\tau''}{r_1 r_2} \mid \mathcal{E}_0\right) \\ & \geq \mathbb{P}\left(z_1 + z_2 \geq \frac{\tau' + \tau''}{\alpha^3} - \frac{\beta^2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#)}{\alpha}\right) - \frac{1}{2} \mathbb{P}\left(|z_3| \geq \frac{\tau''}{\beta^2}\right). \end{aligned}$$

Because  $z_1 + z_2 \sim \mathcal{N}(0, \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp} + \mathbf{P}_{\mathbf{u}_\#} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#}\|_{\mathbb{F}}^2/2)$ , in order to lower-bound the tail probability of  $z_1 + z_2$ , we need to compute a lower estimate of its variance. It follows from (39) that every  $\mathbf{H} \in \mathcal{A}_\delta$  satisfies

$$\|\mathcal{P}_{T^\perp}(\mathbf{H})\|_* \leq \frac{1 - \lambda + \delta}{\lambda}. \quad (\text{H.6})$$

By (H.3) and (H.6) together with Hölder's inequality, we also obtain

$$\|P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}^{\perp}}\|_{\mathbb{F}}^2 \leq \|P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}^{\perp}}\| \cdot \|P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}^{\perp}}\|_* \leq \frac{(1-\lambda+\delta)\zeta}{\lambda}. \quad (\text{H.7})$$

Furthermore, by Lemma 5.1, every  $H \in \mathcal{R}_{\delta}$  satisfies

$$\frac{1-\delta^2}{\delta^2} \cdot |\text{Im}(u_{\sharp}^*Hv_{\sharp})|^2 \leq \|H - P_{u_{\sharp}}HP_{v_{\sharp}}\|_{\mathbb{F}}^2.$$

Then it follows that:

$$\begin{aligned} 1 &= \|H\|_{\mathbb{F}}^2 = \|P_{u_{\sharp}}HP_{v_{\sharp}} + P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}} + P_{u_{\sharp}}HP_{v_{\sharp}^{\perp}}\|_{\mathbb{F}}^2 + |\text{Im}(u_{\sharp}^*Hv_{\sharp})|^2 + |\text{Re}(u_{\sharp}^*Hv_{\sharp})|^2 \\ &\leq \frac{\|P_{u_{\sharp}}HP_{v_{\sharp}} + P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}} + P_{u_{\sharp}}HP_{v_{\sharp}^{\perp}}\|_{\mathbb{F}}^2}{1-\delta^2} + |\text{Re}(u_{\sharp}^*Hv_{\sharp})|^2. \end{aligned} \quad (\text{H.8})$$

It also follows from (39) that every  $H \in \mathcal{A}_{\delta}$  satisfies:

$$\text{Re}(u_{\sharp}^*Hv_{\sharp}) \geq \frac{-\delta}{1-\lambda}. \quad (\text{H.9})$$

The assumption  $\lambda + \delta < 1$  implies that the right-hand side of (H.9) is strictly larger than  $-1$ .

By (H.9) and  $\text{Re}(u_{\sharp}^*Hv_{\sharp}) \leq 0$ , we also have

$$|\text{Re}(u_{\sharp}^*Hv_{\sharp})| \leq \frac{\delta}{1-\lambda}. \quad (\text{H.10})$$

Therefore, by applying (H.10) to (H.8), after a rearrangement, we obtain

$$\|P_{u_{\sharp}}HP_{v_{\sharp}} + P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}} + P_{u_{\sharp}}HP_{v_{\sharp}^{\perp}}\|_{\mathbb{F}}^2 \geq (1-\delta^2) \left(1 - \frac{\delta^2}{(1-\lambda)^2}\right).$$

Then (H.7) implies

$$\|P_{u_{\sharp}}HP_{v_{\sharp}} + P_{u_{\sharp}^{\perp}}HP_{v_{\sharp}}\|_{\mathbb{F}}^2 \geq (1-\delta^2) \left(1 - \frac{\delta^2}{(1-\lambda)^2}\right) - \frac{(1-\lambda+\delta)\zeta}{\lambda}. \quad (\text{H.11})$$

Now, from (H.10) and (H.11), we obtain

$$\mathbb{P}\left(z_1 + z_2 \geq \frac{\tau' + \tau''}{\alpha^3} - \frac{\beta^2 \text{Re}(u_{\sharp}^*Hv_{\sharp})}{\alpha}\right) \geq \mathbb{P}\left(z_1 + z_2 \geq \frac{\tau' + \tau''}{\alpha^3} + \frac{\beta^2\delta}{\alpha(1-\lambda)}\right) \geq \mathbb{P}\left(g \geq \frac{t}{\sigma_{\zeta}}\right) \quad (\text{H.12})$$

for  $g \sim \mathcal{N}(0, 1)$ , where

$$\sigma_{\zeta} = \sqrt{(1-\delta^2) \left(1 - \frac{\delta^2}{(1-\lambda)^2}\right) - \frac{(1-\lambda+\delta)\zeta}{\lambda}}$$

and

$$t = \frac{\tau' + \tau''}{\alpha^3} + \frac{\beta^2\delta}{\alpha(1-\lambda)}. \quad (\text{H.13})$$

Moreover, the tail bound of  $z_3$  in (H.1) implies

$$\begin{aligned} \mathbb{P}\left(|z_3| \geq \frac{\tau''}{\beta^2}\right) &\leq C \exp\left(-\frac{2}{C}\left[\frac{\tau''^2/\beta^4}{\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp}\|_{\text{F}}^2} \wedge \frac{\tau''/\beta^2}{\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp}\|}\right]\right) \\ &\leq C \exp\left(-\frac{2}{C}\left[\frac{\lambda \tau''^2}{\beta^4(1-\lambda+\delta)\zeta} \wedge \frac{\tau''}{\beta^2 \zeta}\right]\right). \end{aligned} \quad (\text{H.14})$$

Note that the tail bound in (H.12) is monotone decreasing in  $t/\sigma_\zeta$ . Furthermore, for those  $\zeta$  that make  $\sigma_\zeta$  positive,  $t/\sigma_\zeta$  is a monotone increasing in  $\zeta$ . (The condition  $\delta \leq 0.2$  implies the existence of such  $\zeta$ .) Hence, the tail bound in (H.12) is monotone decreasing in  $\zeta$ . On the contrary, the upper bound in (H.14) monotonically converges to 0 as  $\zeta > 0$  decreases toward 0. Therefore, there exists small enough  $\zeta$  such that the upper bound in (H.14) becomes less than half of (H.12). Then the lower bound (H.2) is further bounded from below by the half of (H.12). Note that  $\zeta$  is determined independent from all dimension parameters and hence both  $\zeta$  and the resulting lower bound for the probability in (H.2) are numerical constants.

Next we consider the complimentary subcase where  $\text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) > 0$ . Similarly to the previous subcase, since  $z_3$  has a symmetric distribution, it follows that:

$$\begin{aligned} &\mathbb{P}\left(r_1 r_2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{r_1 r_2} \mid \mathcal{E}_0\right) - \mathbb{P}\left(z_3 < -\frac{\tau''}{r_1 r_2} \mid \mathcal{E}_0\right) \\ &\geq \mathbb{P}\left(\alpha^2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{\alpha^2} \mid \mathcal{E}_0\right) - \frac{1}{2} \mathbb{P}\left(|z_3| \geq \frac{\tau''}{\beta^2}\right). \end{aligned}$$

If  $\text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) \leq \alpha^{-4}(\tau' + \tau'')$ , since  $z_1 + z_2$  is a zero-mean Gaussian variable, then it follows that:

$$\mathbb{P}\left(\alpha^2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{\alpha^2} \mid \mathcal{E}_0\right) \geq \mathbb{P}\left(z_1 + z_2 \geq \frac{\tau' + \tau''}{\alpha^3} - \alpha \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#)\right).$$

Thus by choosing  $\tau' + \tau''$  small enough one can satisfy (H.10). Thus, we obtain the desired conclusion as in the previous subcase by repeating the same arguments.

If  $\text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) > \alpha^{-4}(\tau' + \tau'')$  on the other hand, then

$$\mathbb{P}\left(\alpha^2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 \geq \frac{\tau' + \tau''}{\alpha^2} \mid \mathcal{E}_0\right) \geq \mathbb{P}\left(z_1 + z_2 \geq \underbrace{\frac{\tau' + \tau''}{\alpha^2 \beta} - \frac{\alpha^2 \text{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#)}{\beta}}_{<0}\right) > \frac{1}{2},$$

which is larger than the other lower bounds on the tail probability.

**Case 2:** Next we consider the complementary case where

$$\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{H} \mathbf{P}_{\mathbf{v}_\#^\perp}\|_{\text{F}} > \zeta, \quad (\text{H.15})$$

where  $\zeta$  is the constant determined in the previous case. In this case, the lower estimate in (H.2) is further bounded from below by

$$\begin{aligned}
 & \mathbb{P}\left(r_1 r_2 \operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#) + r_2 z_1 + r_1 z_2 + z_3 \geq \frac{\tau'}{r_1 r_2} \mid \mathcal{E}_0\right) \\
 & \geq \mathbb{P}\left(-\beta^2 [\operatorname{Re}(\mathbf{u}_\#^* \mathbf{H} \mathbf{v}_\#)]_- + r_2 z_1 + r_1 z_2 + z_3 \geq \frac{\tau'}{r_1 r_2} \mid \mathcal{E}_0\right) \\
 & \geq \mathbb{P}\left(r_2 z_1 + r_1 z_2 + z_3 \geq \frac{\tau'}{r_1 r_2} + \beta^2(1 - \zeta) \mid \mathcal{E}_0\right) \\
 & \geq \mathbb{P}\left(z_3 \geq \frac{\tau' + \tau''}{r_1 r_2} + \beta^2(1 - \zeta) \mid \mathcal{E}_0\right) + \mathbb{P}\left(r_2 z_1 + r_1 z_2 \geq -\frac{\tau''}{r_1 r_2} \mid \mathcal{E}_0\right) - 1 \\
 & \geq \mathbb{P}\left(z_3 \geq \frac{\tau' + \tau''}{\alpha^2} + \beta^2(1 - \zeta)\right) - \mathbb{P}\left(z_1 + z_2 \geq \frac{\tau''}{\beta^3}\right) \\
 & \geq \frac{1}{2} \mathbb{P}\left(|z_3| \geq \frac{\tau' + \tau''}{\alpha^2} + \beta^2(1 - \zeta)\right) - \mathbb{P}\left(z_1 + z_2 \geq \frac{\tau''}{\beta^3}\right),
 \end{aligned}$$

where the second and third steps follow from (H.10) and the inclusion–exclusion principle, respectively. Then, by (H.15), the tail bound on  $z_3$  is lower-bounded by

$$\mathbb{P}(|z_3| \geq t) \geq \frac{1}{C} \exp\left(-2C\left[\frac{t^2}{\zeta^2} \wedge \frac{t}{\zeta}\right]\right), \quad (\text{H.16})$$

where  $t$  is given in (H.13). Since  $\|\mathbf{H}\|_F = 1$ , the variance of  $z_1 + z_2$  is no larger than  $1/2$ . Thus, the tail bound of  $z_1 + z_2$  is upper-bounded by

$$\mathbb{P}\left(z_1 + z_2 \geq \frac{\tau''}{\beta^3}\right) \leq \exp\left(-\frac{2\tau''^2}{\beta^6}\right). \quad (\text{H.17})$$

Note that  $\tau''$  still remains a free parameter. For every  $t > \zeta$  the lower bound in (H.16) is an exponential tail while the upper bound in (H.17) is a subgaussian tail. Therefore, as  $\tau''$  increases while the other parameters are fixed, by (H.13),  $t$  also increases as an affine function of  $\tau''$  and the lower bound in (H.16) decays slower than the upper bound in (H.17). We may choose  $\tau''$  so that the lower bound in (H.16) is larger than four times the upper bound in (H.17). Then the lower bound (H.2) is further bounded below by the resulting value of (H.17). Again, this lower bound is a numerical constant independent of scaling of all dimension parameters.

### I. Proof of Lemma 5.7

Without loss of generality, we may assume that  $\|\mathbf{H}\|_F = \|\mathbf{X}_\#\|_F = 1$ . Then  $\mathbf{X}_\#$  is written as  $\mathbf{u}_\# \mathbf{v}_\#^*$  where  $\mathbf{u}_\# \in \mathbb{C}^{d_1}$  and  $\mathbf{v}_\# \in \mathbb{C}^{d_2}$  satisfy  $\|\mathbf{u}_\#\|_2 = \|\mathbf{v}_\#\|_2 = 1$ . With this expression of  $\mathbf{X}_\#$ , the Rademacher

complexity  $\mathfrak{C}_M(\mathcal{A}_\delta)$  is written as

$$\begin{aligned}
\mathfrak{C}_M(\mathcal{A}_\delta) &= \mathbb{E} \sup_{\mathbf{H} \in \mathcal{A}_\delta} \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \operatorname{Re}(\mathbf{b}_m^* \mathbf{v}_\# \mathbf{u}_\#^* \mathbf{a}_m \mathbf{a}_m^* \mathbf{H} \mathbf{b}_m) \\
&= \mathbb{E} \sup_{\mathbf{H} \in \mathcal{A}_\delta} \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \operatorname{Re}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^* \mathbf{H}) \\
&\leq \mathbb{E} \sup_{\mathbf{H} \in \mathcal{A}_\delta} \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \operatorname{Re}(\mathcal{P}_T(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*), \mathcal{P}_T(\mathbf{H})) \\
&\quad + \mathbb{E} \sup_{\mathbf{H} \in \mathcal{A}_\delta} \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \operatorname{Re}(\mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*), \mathcal{P}_{T^\perp}(\mathbf{H})) \\
&\leq \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\|_{\mathbb{F}} \cdot \sup_{\mathbf{H} \in \mathcal{A}_\delta} \|\mathcal{P}_T(\mathbf{H})\|_{\mathbb{F}} \\
&\quad + \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\| \cdot \sup_{\mathbf{H} \in \mathcal{A}_\delta} \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_*,
\end{aligned}$$

where the first inequality is obtained by taking the supremum of each summand after applying  $\mathbf{H} = \mathcal{P}_T(\mathbf{H}) + \mathcal{P}_{T^\perp}(\mathbf{H})$  and the second inequality holds by Hölder's inequality.

Since  $\mathcal{P}_T$  is an orthogonal projection onto a subspace, we have  $\|\mathcal{P}_T(\mathbf{H})\|_{\mathbb{F}} \leq \|\mathbf{H}\|_{\mathbb{F}} = 1$ . Furthermore, for all  $\mathbf{H} \in \mathcal{A}_\delta$ ,  $\|\mathcal{P}_{T^\perp}(\mathbf{H})\|_*$  is upper-bounded by (H.6). Therefore, we obtain

$$\begin{aligned}
\mathfrak{C}_M(\mathcal{A}_\delta) &\leq \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\|_{\mathbb{F}} \\
&\quad + \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\| \cdot \left( \frac{1 - \lambda + \delta}{\lambda} \right).
\end{aligned} \tag{I.1}$$

It remains to compute upper estimates of the expectation terms in (I.1). Since  $(\epsilon_m)_{m=1}^M$  is a Rademacher sequence, we have

$$\begin{aligned}
\mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\|_{\mathbb{F}} &\leq \sqrt{\mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_T(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\|_{\mathbb{F}}^2} \\
&= \sqrt{\mathbb{E} \frac{1}{M} \sum_{m=1}^M \|\mathcal{P}_T(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*)\|_{\mathbb{F}}^2} = \sqrt{\mathbb{E} \|\mathcal{P}_T(\mathbf{a} \mathbf{a}^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b} \mathbf{b}^*)\|_{\mathbb{F}}^2},
\end{aligned}$$

where the first step follows from Jensen's inequality and the last step follows since  $\mathbf{a}_1, \dots, \mathbf{a}_M$  (resp.  $\mathbf{b}_1, \dots, \mathbf{b}_M$ ) are independent copies of  $\mathbf{a}$  (resp.  $\mathbf{b}$ ).

Note that  $\mathcal{P}_T(\mathbf{a}\mathbf{a}^*\mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}\mathbf{b}^*)$  is written as

$$\begin{aligned} \mathcal{P}_T(\mathbf{a}\mathbf{a}^*\mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}\mathbf{b}^*) &= \mathbf{a}^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b} \cdot \mathcal{P}_T(\mathbf{a}\mathbf{b}^*) \\ &= \mathbf{a}^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b} \cdot (\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#} + \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#} + \mathbf{P}_{\mathbf{u}_\#} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#^\perp}), \end{aligned}$$

where  $\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#}$ ,  $\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#}$  and  $\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#^\perp}$  are mutually orthogonal matrices in the Hilbert space  $S_2$ . Thus, the Pythagorean identity implies

$$\begin{aligned} \|\mathcal{P}_T(\mathbf{a}\mathbf{a}^*\mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}\mathbf{b}^*)\|_{\mathbb{F}}^2 &= |\mathbf{a}^* \mathbf{u}_\#|^2 |\mathbf{b}^* \mathbf{v}_\#|^2 (\|\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#}\|_{\mathbb{F}}^2 + \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#}\|_{\mathbb{F}}^2 + \|\mathbf{P}_{\mathbf{u}_\#} \mathbf{a}\mathbf{b}^* \mathbf{P}_{\mathbf{v}_\#^\perp}\|_{\mathbb{F}}^2) \\ &= |\mathbf{a}^* \mathbf{u}_\#|^4 |\mathbf{b}^* \mathbf{v}_\#|^4 + |\mathbf{a}^* \mathbf{u}_\#|^2 |\mathbf{b}^* \mathbf{v}_\#|^4 \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}\|_2^2 + |\mathbf{a}^* \mathbf{u}_\#|^4 |\mathbf{b}^* \mathbf{v}_\#|^2 \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}\|_2^2. \end{aligned}$$

Since  $\mathbf{a} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{d_1})$  and  $\mathbf{b} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_{d_2})$  are independent,  $\mathbf{a}^* \mathbf{u}_\#$ ,  $\mathbf{b}^* \mathbf{v}_\#$ ,  $\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}$  and  $\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}$  are all mutually independent. Therefore, exploiting this independence, one can show that the expectation is upper-bounded by

$$\mathbb{E} \|\mathcal{P}_T(\mathbf{a}\mathbf{a}^*\mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}\mathbf{b}^*)\|_{\mathbb{F}}^2 \leq 2\|\mathbf{u}_\#\|_2^2 \|\mathbf{v}_\#\|_2^2 (2 + d_1 + d_2).$$

By Jensen's inequality, the second expectation in (I.1) is upper-bounded by

$$\mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\| \leq \left( \mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\|^p \right)^{1/p} \quad (\text{I.2})$$

for all  $p \in 2\mathbb{N}$ . To upper bound the right-hand side of (I.2), we apply Theorem B.2 for

$$\mathbf{Y}_m = \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) = \epsilon_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}_m \mathbf{b}_m^* \mathbf{P}_{\mathbf{v}_\#^\perp}, \quad m = 1, \dots, M,$$

with some  $p \in \mathbb{N}$  that satisfies  $p \geq 2$ . Note that  $\mathbb{E} \mathbf{Y}_m = \mathbf{0}$  for all  $m = 1, \dots, M$ . By direct computation, we obtain

$$\mathbb{E} \mathbf{Y}_m \mathbf{Y}_m^* = \|\mathbf{u}_\#\|_2^2 \|\mathbf{v}_\#\|_2^2 \text{tr}(\mathbf{P}_{\mathbf{v}_\#^\perp}) \mathbf{P}_{\mathbf{u}_\#^\perp} \quad \text{and} \quad \mathbb{E} \mathbf{Y}_m^* \mathbf{Y}_m = \|\mathbf{u}_\#\|_2^2 \|\mathbf{v}_\#\|_2^2 \text{tr}(\mathbf{P}_{\mathbf{u}_\#^\perp}) \mathbf{P}_{\mathbf{v}_\#^\perp}, \quad m = 1, \dots, M.$$

Therefore,

$$\left\| \sum_{m=1}^M \mathbb{E} \mathbf{Y}_m \mathbf{Y}_m^* \right\|^{1/2} \vee \left\| \sum_{m=1}^M \mathbb{E} \mathbf{Y}_m^* \mathbf{Y}_m \right\|^{1/2} \leq \|\mathbf{u}_\#\|_2 \|\mathbf{v}_\#\|_2 \sqrt{M(d_1 + d_2)}.$$

Since the spectral norm of  $\mathbf{Y}_m$  is upper-bounded by

$$\|\mathbf{Y}_m\| = |\mathbf{a}_m^* \mathbf{u}_\#| |\mathbf{b}_m^* \mathbf{v}_\#| \|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}_m\|_2 \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}_m\|_2 \leq 2|\mathbf{a}_m^* \mathbf{u}_\#| |\mathbf{b}_m^* \mathbf{v}_\#| (\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}_m\|_2^2 + \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}_m\|_2^2),$$

it follows that:

$$\begin{aligned} (\mathbb{E} \|\mathbf{Y}_m\|^p)^{1/p} &\leq 2(\mathbb{E} |\mathbf{a}_m^* \mathbf{u}_\#|^p)^{1/p} \cdot (\mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^p)^{1/p} \cdot [\mathbb{E} (\|\mathbf{P}_{\mathbf{u}_\#^\perp} \mathbf{a}_m\|_2^2 + \|\mathbf{P}_{\mathbf{v}_\#^\perp} \mathbf{b}_m\|_2^2)^p]^{1/p} \\ &\leq 2(\mathbb{E} |\mathbf{a}_m^* \mathbf{u}_\#|^p)^{1/p} \cdot (\mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^p)^{1/p} \cdot [\mathbb{E} (\|\mathbf{a}_m\|_2^2 + \|\mathbf{b}_m\|_2^2)^p]^{1/p}. \end{aligned}$$

Since  $\mathbf{a}_m^* \mathbf{u}_\# \sim \mathcal{CN}(0, 1)$  and  $\mathbf{b}_m^* \mathbf{v}_\# \sim \mathcal{CN}(0, 1)$ , we have

$$(\mathbb{E} |\mathbf{a}_m^* \mathbf{u}_\#|^p)^{1/p} = (\mathbb{E} |\mathbf{b}_m^* \mathbf{v}_\#|^p)^{1/p} \leq C_1 \sqrt{p}.$$

for a numerical constant  $C_1$ . Since  $2(\|\mathbf{a}_m\|_2^2 + \|\mathbf{b}_m\|_2^2)$  is a chi-square random variable of the degree-of-freedom  $2(d_1 + d_2)$ , it follows that for  $p \geq 2$  we have:

$$(\mathbb{E}(\|\mathbf{a}_m\|_2^2 + \|\mathbf{b}_m\|_2^2)^p)^{1/p} \leq 2d_1 + 2d_2 + C_2 p$$

for a numerical constant  $C_2$ . By collecting these estimates, we obtain

$$p \left( \sum_{m=1}^M \mathbb{E} \|\mathbf{Y}_m\|^p \right)^{1/p} \leq C_3 M^{1/p} p^2 (d_1 + d_2 + p).$$

Applying the above estimates to Theorem B.2 together with (I.2) provides

$$\mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\| \leq C_4 \left( \sqrt{p(d_1 + d_2)} + M^{1/p-1/2} p^2 (d_1 + d_2 + p) \right). \quad (\text{I.3})$$

As we set  $p = \log M$ , (I.3) implies

$$\mathbb{E} \left\| \frac{1}{\sqrt{M}} \sum_{m=1}^M \epsilon_m \mathcal{P}_{T^\perp}(\mathbf{a}_m \mathbf{a}_m^* \mathbf{u}_\# \mathbf{v}_\#^* \mathbf{b}_m \mathbf{b}_m^*) \right\| \leq C_5 \sqrt{(d_1 + d_2) \log M} \cdot \left( 1 + \frac{\sqrt{d_1 + d_2} \log^{3/2} M}{\sqrt{M}} \right).$$

Then (10) implies that the right-hand side is further upper bounded by  $C_6 \sqrt{d_1 + d_2} \log M$ . Finally, (42) is obtained by plugging in these upper estimates to (I.1).