# Finite Sample Guarantees for PCA in Non-Isotropic and Data-Dependent Noise

Namrata Vaswani and Praneeth Narayanamurthy
{`namrata, pkurpadn`}@iastate.edu
Department of Electrical and Computer Engineering
Iowa State University

## Abstract

This work obtains novel finite sample guarantees for Principal Component Analysis (PCA). These hold even when the corrupting noise is non-isotropic, and a part (or all of it) is data-dependent. Because of the latter, in general, the noise and the true data are correlated. The results in this work are a significant improvement over those given in our earlier work where this "correlated-PCA" problem was first studied. In fact, in certain regimes, our results imply that the sample complexity required to achieve subspace recovery error that is a constant fraction of the noise level is near-optimal. Useful corollaries of our result include guarantees for PCA in sparse data-dependent noise and for PCA with missing data. An important application of the former is in proving correctness of the subspace update step of a popular online algorithm for dynamic robust PCA.

## 1  Introduction

Principal Components Analysis (PCA) is among the most frequently used tools for dimension reduction for a wide variety of data analysis applications. Some examples include exploratory data analysis, data classification, image or video retrieval, face recognition, and recommendation system design. Given a matrix of observed data, the goal of PCA is to compute a small number of orthogonal directions that contain most of the variability of the data. These principal components are easily computed via singular value decomposition (SVD) on the observed data matrix.

PCA is a classical and very well-studied problem. There has been a large amount of work on analyzing PCA, however most existing results for PCA are asymptotic, e.g., see [1], and references therein. While asymptotic analysis is useful because it provides limits on what can be done, it is less practically relevant. A very nice work by Nadler [2] provides finite sample guarantees for one-dimensional PCA that hold under the spiked covariance model [3]. Spiked covariance model means that true data ("signal") and noise are independent,

1

or at least uncorrelated, and the noise is isotropic (noise power in all directions is equal). A simple example of isotropic noise is noise that is zero mean with a covariance matrix that is a scalar multiple of identity. There is also much new work on analyzing streaming solutions for PCA in the non-asymptotic setting, e.g. [4, 5], however that is a different problem (places memory constraints on the algorithm) and we will not discuss it here. All of the above works either assume the spiked covariance model [3, 1, 2, 4] or only analyze one-dimensional PCA [5, 1, 2] or both [1, 2].

Our work obtains novel finite sample (non-asymptotic) guarantees for $r$-dimensional PCA (with $r \geq 1$) that hold even when the corrupting noise is non-isotropic, and, a part, or all, of it is data-dependent. Because of the latter, in general, the data and noise are no longer independent (or even uncorrelated). A special case of this problem was first studied in [6] where we called it "correlated-PCA". As we will explain, the current work significantly improves upon the results of [6].

*Notation.* We use $\boldsymbol{A}'$ to denote transpose of a matrix $\boldsymbol{A}$. We use $\|\cdot\|_p$ to denote the $l_p$ norm of a vector or the induced $l_p$ norm of a matrix. Most of this paper only uses $l_2$ norm. At a few places, even if the subscript is missing, it refers to the $l_2$ norm. For a set of indices $\mathcal{T}$, $\boldsymbol{I}_{\mathcal{T}}$ refers to an $n \times |\mathcal{T}|$ matrix of columns of the identity matrix indexed by entries in $\mathcal{T}$. For a matrix $\boldsymbol{A}$, $\boldsymbol{A}_{\mathcal{T}} := \boldsymbol{A}\boldsymbol{I}_{\mathcal{T}}$. A tall matrix, $\boldsymbol{P}$, with orthonormal columns is referred to as a *basis matrix*. We use span($\boldsymbol{P}$) to denote the span of the columns of the basis matrix $\boldsymbol{P}$ and we use $\boldsymbol{P}_{\perp}$ to denote a basis matrix whose span is the orthogonal complement of span($\boldsymbol{P}$). Thus $\boldsymbol{P}\boldsymbol{P}' + \boldsymbol{P}_{\perp}\boldsymbol{P}_{\perp}' = \boldsymbol{I}$. For two *basis matrices* $\hat{\boldsymbol{P}}$, $\boldsymbol{P}$, we define the subspace recovery error (SE) as

$$\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) := \|(\boldsymbol{I} - \hat{\boldsymbol{P}}\hat{\boldsymbol{P}}')\boldsymbol{P}\|_2.$$

This measures the sine of the principal angle between column spans of $\hat{\boldsymbol{P}}$ and $\boldsymbol{P}$.

We re-use the letters $C, c$ at various places to denote different numerical constants. Also, $\frac{1}{\alpha}\sum_t f(t)$ is often used instead of $\frac{1}{\alpha}\sum_{t=1}^{\alpha} f(t)$.

*Problem Setting.* We study PCA in the following setting which assumes that the data-dependent component of the noise at each time $t$ depends linearly on the true data (signal) vector at time $t$. For $t = 1, 2, \ldots, \alpha$, we are given $n$-length observed data vectors, $\boldsymbol{y}_t$, that satisfy

$$\boldsymbol{y}_t := \boldsymbol{\ell}_t + \boldsymbol{w}_t + \boldsymbol{v}_t, \text{ where } \boldsymbol{\ell}_t = \boldsymbol{P}\boldsymbol{a}_t, \ \boldsymbol{w}_t = \boldsymbol{M}_t\boldsymbol{\ell}_t, \ \mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{v}_t'] = 0, \qquad (1)$$

$\boldsymbol{P}$ is an $n \times r$ basis matrix with $r \ll n$; $\boldsymbol{\ell}_t$ is the true data ("signal") vector that lies in an $r$ dimensional subspace of $\mathbb{R}^n$, span($\boldsymbol{P}$); $\boldsymbol{a}_t$ is its projection into this subspace; $\boldsymbol{w}_t$ is the data-dependent noise component; and $\boldsymbol{v}_t$ is the uncorrelated noise component, i.e., it satisfies $\mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{v}_t'] = 0$. The data-dependency matrices $\boldsymbol{M}_t$ are *unknown* and such that the signal-noise correlation $\mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{w}_t'] \neq 0$. Thus, we also often refer to $\boldsymbol{w}_t$ as "correlated" noise. The goal is to estimate span($\boldsymbol{P}$). Since the matrices $\boldsymbol{M}_t$ are *time-varying*, observe that, the $\boldsymbol{w}_t$'s taken together, in general, do not lie in a lower dimensional subspace of $\mathbb{R}^n$.

Data-dependent noise occurs in a large number of applications due to signal reflections or signal leakage, e.g., in electro-encephalography (EEG) and magneto-encephalography (MEG). It is called interference in these settings. It also often occurs in molecular biology applications when the noise affects the measurement levels through the very same process as the interesting signal [7]. Two other examples where it occurs include PCA with missing data and the subspace update step of the Recursive Projected Compressive Sensing (ReProCS) solution to dynamic robust PCA [8, 9]. In these last two examples, the noise also satisfies our required assumption on signal-noise correlation. We explain them in detail in Sec. 4.1 and 4.2.

Non-isotropic noise is even more common. In signal processing literature, it is often referred to as "colored" noise. One common example of this is the noise in different pixels of an image sequence [10]. The variance of the noise is often region-dependent and this is what results in the non-isotropy.

*The SVD solution.* This computes $\hat{\boldsymbol{P}}$ as the top $r$ left singular vectors of the observed data matrix $[\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_\alpha]$. Equivalently $\hat{\boldsymbol{P}}$ is the matrix of top $r$ eigenvectors of the $n \times n$ matrix $\boldsymbol{D} := \frac{1}{\alpha} \sum_{t=1}^{\alpha} \boldsymbol{y}_t \boldsymbol{y}_t'$. Hence this solution is often also referred to as *EVD (eigenvalue decomposition)*.

*Related Work.* To our best knowledge, existing guarantees for PCA other than [2, 6] are asymptotic. Also, see the discussion in [4, Section 1]. We discuss these and two other tangentially related works [11, 12] in Sec. 3.

*Contributions.* This work, which builds on work in [6], is the first to study PCA in a non-isotropic and data-dependent noise setting. Our main result (Theorem 2.13) shows that it is possible to recover the signal subspace, $\text{span}(\boldsymbol{P})$, with error at most $\varepsilon$ as long as (a) a simple assumption on signal-noise correlation holds, (b) the ratio between the maximum signal-noise correlation and the minimum signal subspace eigenvalue is upper bounded; (c) the ratio between the noise power outside the signal subspace and the minimum signal subspace eigenvalue is upper bounded; and (d) the sample complexity, $\alpha$, is lower bounded. All the required bounds depend on $\varepsilon$. We obtain such a result in two settings - bounded signal and noise and sub-Gaussian (e.g., Gaussian) signal and noise. In most applications, boundedness is a more practical assumption than Gaussianity since data acquisition devices usually have bounded power.

As compared to the result of [6], our results holds under a much weaker signal-noise correlation assumption *and* needs a sample complexity lower bound that is much better than the one given in [6]. In fact, for the only data-dependent noise case studied in [6], our sample complexity bound is near-optimal. Secondly, we generalize the observed data model to also include an uncorrelated, but possibly non-isotropic, noise term. This is a more practically valid noise model since the noise/corruption is usually not fully data-dependent. Moreover, this allow us to obtain the existing isotropic noise results as special cases. Lastly, we also provide a simple provably correct method for automatic subspace dimension estimation that does not use knowledge of any model parameter (see Theorem 2.16).

*Paper Organization.* We state and discuss the main results in Sec. 2. Related works are discussed in Sec. 3. In Sec. 4, we show how our result can be

applied to the problem of PCA in sparse data-dependent noise (PCA-SDDN) and to its two special cases - PCA in missing data, and the subspace update step of ReProCS for dynamic robust PCA. Numerical experiments backing our theoretical claims are shown in Sec. 5. We conclude in Sec. 6.

# 2 Main Results

In Sec. 2.1, we state our basic assumptions and define a few quantities. In Sec. 2.2 and 2.3, we state and discuss corollaries for the two special cases - data and only uncorrelated noise and data and only data-dependent noise. We give the most general version of our result (Theorem 2.13) in Sec. 2.4. This and its corollaries assume that the subspace dimension, $r$, is known. We show how to provably correct estimate $r$ (Theorems 2.15 and 2.16) in Sec. 2.5.

## 2.1 Basic assumptions

In this entire paper, we assume the following.

**Asssumption 2.1.** *The $\boldsymbol{\ell}_t$'s satisfy $\boldsymbol{\ell}_t = \boldsymbol{P}\boldsymbol{a}_t$ with $\boldsymbol{a}_t$'s being zero mean and mutually independent random variables (r.v.), with diagonal covariance matrix, $\boldsymbol{\Lambda} := \mathbb{E}[\boldsymbol{a}_t\boldsymbol{a}_t']$.*

*The $\boldsymbol{v}_t$'s are zero mean and mutually independent r.v.'s, with covariance matrix $\boldsymbol{\Sigma}_v := \mathbb{E}[\boldsymbol{v}_t\boldsymbol{v}_t']$. Also, $\mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{v}_t'] = 0$ for all $t$, i.e., they are uncorrelated.*

*(Notice that the model on $\boldsymbol{\ell}_t$ automatically imposes a model on the data-dependent noise component $\boldsymbol{w}_t := \boldsymbol{M}_t\boldsymbol{\ell}_t$.)*

We define a few quantities to state our results compactly.

**Definition 2.2.** *Let*

1. $\lambda^- := \lambda_{\min}(\boldsymbol{\Lambda})$, $\lambda^+ := \lambda_{\max}(\boldsymbol{\Lambda})$ *and*

$$f := \frac{\lambda^+}{\lambda^-}.$$

2. *Define the following functions of $\boldsymbol{\Sigma}_v$:*

$$\lambda_{v,\boldsymbol{P}}^- := \lambda_{\min}(\boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}), \ \lambda_{v,\mathrm{rest}}^+ := \lambda_{\max}(\boldsymbol{\Sigma}_v - \boldsymbol{P}\boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}\boldsymbol{P}'), \ \lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp} := \|\boldsymbol{P}_\perp'\boldsymbol{\Sigma}_v\boldsymbol{P}\|_2$$

*and*

$$\lambda_v^+ := \|\boldsymbol{\Sigma}_v\|_2.$$

*It is easy to see that $\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp} \leq \lambda_{v,\mathrm{rest}}^+$. Also, $\lambda_{v,\mathrm{rest}}^+ \leq \lambda_v^+$, $\lambda_{v,\boldsymbol{P}}^- \leq \lambda_v^+$.*

3. *The following factor will used at various places in our results:*

$$g := \max\left(\frac{\lambda_v^+}{\lambda^-}, \sqrt{\frac{\lambda_v^+}{\lambda^-}f}\right).$$

We assume that $\lambda_v^+$ and $\lambda^+$ are at most constant ($O(1)$) with $n$.

**Remark 2.3.** *For notational simplicity, we have let $\boldsymbol{\Lambda}$ and $\boldsymbol{\Sigma}_v$ be constant with $t$. However, all our proofs will go through with minor changes if these are time-varying. We explain the changes needed in Remark 2.14.*

**Asssumption 2.4** (Bounded signal and noise). *Assumption 2.1 holds and the $\boldsymbol{a}_t$'s are element-wise bounded r.v.'s, i.e., there exists a numerical constant, $\eta$, such that,*

$$\max_{j=1,2,\ldots r} \max_t \frac{(\boldsymbol{a}_t)_j^2}{\lambda_j(\boldsymbol{\Lambda})} \leq \eta.$$

*For example, if $\boldsymbol{a}_t$'s are uniformly distributed, then $\eta = 3$. Throughout this paper, $\eta$ will be treated as a numerical constant.*

*The $\boldsymbol{v}_t$'s are bounded r.v.'s, i.e., there exists an integer $r_v \leq Cn$ such that*

$$\max_t \|\boldsymbol{v}_t\|_2^2 \leq r_v \lambda_v^+.$$

*Here $r_v$ can be interpreted as the "effective noise dimension" of $\boldsymbol{v}_t$.*

**Asssumption 2.5** (Sub-Gaussian signal and noise). *Assumption 2.1 holds and $\boldsymbol{a}_t$'s and $\boldsymbol{v}_t$'s are sub-Gaussian r.v.'s with sub-Gaussian norms bounded by $C\sqrt{\lambda^+}$ and $C\sqrt{\lambda_v^+}$ respectively. Recall from Definition 2.2 that both these are assumed to be $O(1)$ w.r.t. $n$ and so $\boldsymbol{a}_t$'s and $\boldsymbol{v}_t$'s are "nice" sub-Gaussians [13].*

**Remark 2.6.** *We should point out that Assumption 2.4 is not always a special case of Assumption 2.5 even though all bounded r.v.'s are formally sub-Gaussian. It is a special case if we assume that $\boldsymbol{v}_t$ is also element-wise bounded by a constant (w.r.t. $n$). Without this, when $r_v = n$, it is possible that the sub-Gaussian norm of $\boldsymbol{v}_t$ is as large as $\sqrt{n}$ [13]. One example for which this happens is the coordinate distribution [13, Example 5.25]: $\boldsymbol{v}_t$ is equally likely to take one of $2n$ possible values $\{\pm\sqrt{n}\boldsymbol{e}_i\}_{i=1,2,\ldots,n}$ where $\boldsymbol{e}_i$ is the $i$-th column $\boldsymbol{I}$.*

## 2.2 Result for only uncorrelated noise case

Before stating the most general result, we state its corollaries for only uncorrelated and only correlated noise. Also, for simplicity, the results in this and the next subsection assume that the subspace dimension $r$ is known. We explain how to provably correctly estimate $r$ automatically in Sec. 2.5.

**Corollary 2.7** (uncorrelated non-isotropic noise). *Given data vectors $\boldsymbol{y}_t := \boldsymbol{\ell}_t + \boldsymbol{v}_t$ for $t = 1, 2, \ldots, \alpha$ with $\boldsymbol{\ell}_t, \boldsymbol{v}_t$ uncorrelated. Let $\hat{\boldsymbol{P}}$ denote the matrix of top $r$ eigenvectors of $\boldsymbol{D} := \frac{1}{\alpha}\sum_t \boldsymbol{y}_t\boldsymbol{y}_t'$ and define*

$$d(\alpha) := cg\sqrt{\frac{\max(r_v, r)\log n}{\alpha}} \text{ and } d_{\mathrm{denom}}(\alpha) := cf\sqrt{\frac{r + \log n}{\alpha}}.$$

1. *If Assumption 2.4 holds, $\alpha^3 > \max(r_v, r)\log n$, and $\frac{\lambda^+_{v,\mathrm{rest}} - \lambda^-_{v,\boldsymbol{P}}}{\lambda^-} + d(\alpha) + d_{\mathrm{denom}}(\alpha) < 1$, then, w.p. at least $1 - 10n^{-10}$,*

$$\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq \frac{\frac{\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}}{\lambda^-} + d(\alpha)}{1 - \frac{\lambda^+_{v,\mathrm{rest}} - \lambda^-_{v,\boldsymbol{P}}}{\lambda^-} - d(\alpha) - d_{\mathrm{denom}}(\alpha)}.$$

2. *If Assumption 2.5 holds (instead of Assumption 2.4), then we have the same result as above but with $d(\alpha) = d(\alpha)_{sG} := c\max\left(f, \frac{\lambda^+_v}{\lambda^-}\right)\sqrt{\frac{n}{\alpha}}$. Also the result now holds w.p. greater than $1 - 10\exp(-cn)$ and we do not need $\alpha^3 \geq r\log n$.*

*Proof Outline.* This is a corollary of Theorem 2.13 given later and proved in the Appendix. We give the main proof idea here. It relies on the Davis-Kahan $\sin\theta$ theorem [14] which states the following.

**Lemma 2.8** (Davis-Kahan $\sin\theta$ theorem). *Let $\boldsymbol{D}_0$ be a Hermitian matrix whose span of top $r$ eigenvectors equals $\mathrm{span}(\boldsymbol{P})$. Let $\boldsymbol{D}$ be the Hermitian matrix with top $r$ eigenvectors $\hat{\boldsymbol{P}}$. Then,*

$$\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq \frac{\|(\boldsymbol{D} - \boldsymbol{D}_0)\boldsymbol{P}\|_2}{\lambda_r(\boldsymbol{D}_0) - \lambda_{r+1}(\boldsymbol{D})} \leq \frac{\|(\boldsymbol{D} - \boldsymbol{D}_0)\boldsymbol{P}\|_2}{\lambda_r(\boldsymbol{D}_0) - \lambda_{r+1}(\boldsymbol{D}_0) - \lambda_{\max}(\boldsymbol{D} - \boldsymbol{D}_0)} \quad (2)$$

*as long as the denominator is positive. The second inequality follows from the first using Weyl's inequality.*

We apply the above result with $\boldsymbol{D}_0 = \boldsymbol{P}(\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}'_t + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P})\boldsymbol{P}'$ and use $\lambda_r(\boldsymbol{D}_0) \geq \lambda^-_{v,\boldsymbol{P}} + \lambda^- - \|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}'_t - \boldsymbol{\Lambda}\|_2$, $\lambda_{r+1}(\boldsymbol{D}_0) = 0$, and $\boldsymbol{D} - \boldsymbol{D}_0 = \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0] + (\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0])$ to simplify the bound. We then use appropriate concentration inequalities to upper bound $\|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2$ and $\|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}'_t - \boldsymbol{\Lambda}\|_2$. Matrix Bernstein [15] is used for the former and Vershynin's sub-Gaussian result [13, Theorem 5.39] is used for the latter. Since the latter involves $r \times r$ matrices, this gives a better bound - $d_{\mathrm{denom}}(\alpha)$ defined above - than matrix Bernstein would give for this term. $\boxtimes$

A further corollary of the above result essentially recovers the subspace error bound given in [2] for the case of isotropic independent noise (spiked covariance model). This follows because, when $\boldsymbol{\Sigma}_v = \lambda^+_v \boldsymbol{I}$, $\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp} = 0$ and $\lambda^+_{v,\mathrm{rest}} = \lambda^-_{v,\boldsymbol{P}} = \lambda^+_v$. The result of [2] also assumed $r = 1$ and Gaussianity.

**Corollary 2.9** (uncorrelated isotropic noise). *In the setting of Corollary 2.7, if $\boldsymbol{\Sigma}_v = \lambda^+_v \boldsymbol{I}$, then the following simpler result holds: if $\alpha^3 > \max(r_v, r)\log n$, and $d(\alpha) + d_{\mathrm{denom}}(\alpha) < 0.95$, then, w.p. at least $1 - 10n^{-10}$, $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq \frac{d(\alpha)}{1 - d(\alpha) - d_{\mathrm{denom}}(\alpha)}$. In the sub-Gaussian case, $d(\alpha) \equiv d(\alpha)_{sG}$.*

From Corollary 2.9, it is clear that, in case of isotropic noise, to achieve subspace error below $\epsilon$ we only need a lower bound on sample complexity $\alpha$.

6

However, in the general case (Corollary 2.7), we need this sample complexity bound *and* an extra assumption such as the following:

$$\lambda_{v,\text{rest}}^+ - \lambda_{v,\boldsymbol{P}}^- < 0.5\lambda^-, \text{ and } 2\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp} < 0.5\epsilon\lambda^-. \tag{3}$$

To understand why more assumptions are needed in the general case, observe that $\mathbb{E}[\boldsymbol{D}] = \boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}' + \boldsymbol{\Sigma}_v$. If $\boldsymbol{\Sigma}_v = c\boldsymbol{I}$ (isotropic noise), then the span of top $r$ eigenvectors of $\mathbb{E}[\boldsymbol{D}]$ is equal to span($\boldsymbol{P}$). Thus, as long as $\alpha$ is large enough (sample complexity bound holds), by the $\sin\theta$ theorem stated above, the same will be approximately true for span($\hat{\boldsymbol{P}}$) which is the span of top $r$ eigenvectors of $\boldsymbol{D}$. However, when the noise is not isotropic, this is no longer the case. Without extra assumptions, the span of top eigenvectors of $\mathbb{E}[\boldsymbol{D}]$ can be very different from span($\boldsymbol{P}$). We give a simple example below.

**Example 2.10.** *Suppose that $\boldsymbol{\Sigma}_v = (1.2\lambda^-)(\boldsymbol{P}_\perp)_1(\boldsymbol{P}_\perp)_1'$ where $(\boldsymbol{P}_\perp)_1$ is any one direction from* span($\boldsymbol{P}_\perp$); *thus $\boldsymbol{P}'(\boldsymbol{P}_\perp)_1 = 0$. With this,*

$$\mathbb{E}[\boldsymbol{D}] = [\boldsymbol{P}\ (\boldsymbol{P}_\perp)_1] \begin{bmatrix} \boldsymbol{\Lambda} & \\ & 1.2\lambda^- \end{bmatrix} \begin{bmatrix} \boldsymbol{P}' \\ (\boldsymbol{P}_\perp)_1' \end{bmatrix}$$

*(in the above expression, eigenvalues are not in decreasing order). Since $1.2\lambda^- > \lambda^-$, it is clear that the top $r$ eigenvectors of $\mathbb{E}[\boldsymbol{D}]$ will be $[\boldsymbol{P}_1, \boldsymbol{P}_2, \ldots, \boldsymbol{P}_{r-1}, (\boldsymbol{P}_\perp)_1]$ (this statement assumes $\lambda_{r-1}(\boldsymbol{\Lambda}) > \lambda^-$). Thus their span will be orthogonal to $\boldsymbol{P}_r$. As a result the SE between this span and* span($\boldsymbol{P}$) *will be one. Hence when $\alpha$ is large enough so that $\|\boldsymbol{D} - \mathbb{E}[\boldsymbol{D}]\|_2$ is small with high probability (whp), then* SE($\hat{\boldsymbol{P}}, \boldsymbol{P}$) *will also be close to one. To be precise, the following can be shown.*

*Suppose that $\|\boldsymbol{D} - \mathbb{E}[\boldsymbol{D}]\| \leq \epsilon\lambda^-$ for any $\epsilon < 0.01$ and that $\lambda_{r-1}(\boldsymbol{\Lambda}) \geq 1.1\lambda^-$. Then* SE($\hat{\boldsymbol{P}}, \boldsymbol{P}$) $\geq 1 - 11.1\epsilon$.

*To see why this holds, let $\boldsymbol{P}_{\mathbb{E}D} := [\boldsymbol{P}_1, \boldsymbol{P}_2, \ldots, \boldsymbol{P}_{r-1}, (\boldsymbol{P}_\perp)_1]$ denote the top $r$ eigenvectors of $\mathbb{E}[\boldsymbol{D}]$. Notice that $\lambda_r(\mathbb{E}[\boldsymbol{D}]) = \max(1.2\lambda^-, \lambda_{r-1}(\boldsymbol{\Lambda})) \geq 1.1\lambda^-$ and $\lambda_{r+1}(\mathbb{E}[\boldsymbol{D}]) = \lambda^-$. Using Davis-Kahan,*

$$\text{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}_{\mathbb{E}D}) \leq \frac{\|\boldsymbol{D} - \mathbb{E}[\boldsymbol{D}]\|_2}{\lambda_r(\mathbb{E}[\boldsymbol{D}]) - \lambda_{r+1}(\mathbb{E}[\boldsymbol{D}]) - \|\boldsymbol{D} - \mathbb{E}[\boldsymbol{D}]\|_2} \leq \frac{\epsilon\lambda^-}{1.1\lambda^- - \lambda^- - \epsilon\lambda^-} \leq 11.1\epsilon$$

*Using this, triangle inequality, $\boldsymbol{P}_{\mathbb{E}D}'\boldsymbol{P}_r = 0$ and* SE($\boldsymbol{P}_{\mathbb{E}D}, \hat{\boldsymbol{P}}$) = SE($\hat{\boldsymbol{P}}, \boldsymbol{P}_{\mathbb{E}D}$) *(this holds since both $\hat{\boldsymbol{P}}$ and $\boldsymbol{P}_{\mathbb{E}D}$ have the same dimension),*

$$\text{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \geq \text{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}_r) \geq 1 - \|\hat{\boldsymbol{P}}'\boldsymbol{P}_r\|_2 = 1 - \|\hat{\boldsymbol{P}}'(\boldsymbol{I} - \boldsymbol{P}_{\mathbb{E}D}\boldsymbol{P}_{\mathbb{E}D}')\boldsymbol{P}_r\|_2$$
$$\geq 1 - \text{SE}(\boldsymbol{P}_{\mathbb{E}D}, \hat{\boldsymbol{P}}) = 1 - \text{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}_{\mathbb{E}D}) \geq 1 - 11.1\epsilon.$$

For the above example, $\lambda_{v,\boldsymbol{P}}^- = 0$ while $\lambda_{v,\text{rest}}^+ = 1.2\lambda^-$ and hence (3) does not hold. Because of this, the expected value of the average energy of $\boldsymbol{y}_t$'s in a direction outside span($\boldsymbol{P}$) is larger than that in a direction that is in span($\boldsymbol{P}$) and this is what causes SE($\hat{\boldsymbol{P}}, \boldsymbol{P}$) to be large. Assuming (3) helps ensure that

the above does not happen. It also ensures that the maximal correlation between a component of the projection of $\boldsymbol{y}_t$'s in span($\boldsymbol{P}$) and that in span($\boldsymbol{P}_\perp$) is small (bound on $\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}$).

*Sample complexity.* Consider the required lower bound on $\alpha$ to achieve error $\epsilon$. In the bounded case, our result needs $\alpha \geq C \max\left(\frac{g^2}{\epsilon^2}\max(r_v,r)\log n, f^2(r+\log n)\right)$. In the sub-Gaussian case, it needs $\alpha \geq C\frac{g^2}{\epsilon^2}n$. Since the subspace dimension is $r$, the minimum number of samples required in any setting is $r$ (this number suffices only when there is no noise outside span($\boldsymbol{P}$) and all observations are linearly independent). Thus, if $r_v$ is small, e.g., if $r_v \in O(r)$, the sample complexity $\alpha$ required to achieve error $\epsilon$ that is a constant fraction of $g$ is only $(\log n)$ times more, i.e., it is nearly optimal.

On the other hand, $r_v$ can be as larger as $n$. If $r_v = Cn$, and if $\boldsymbol{v}_t$ is also *element-wise* bounded, then it is a "nice" sub-Gaussian and we can use the sub-Gaussian case result to conclude that, in this case, the required sample complexity is $Cn$ (and not $Cn\log n$ as predicted by the bounded case result). However, if $r_v = Cn$ and no other assumption is placed on $\boldsymbol{v}_t$, then, it is not guaranteed to be a "nice" sub-Gaussian (see Remark 2.6). In this case, the required sample complexity will indeed be $C(n\log n)$. We discuss this point further in Sec. 3, where we connect it to the results of [12].

**Remark 2.11.** *It may be possible to reduce the required sample complexity lower bound for the sub-Gaussian case in settings where $r_v \ll n$, e.g., if $r_v = Cr$, or in the data-dependent noise case discussed below in Sec. 2.3. For unbounded noise, we can define $r_v$ as $\mathbb{E}[\|\boldsymbol{v}_t\|_2^2]/\lambda_v^+$. In our current proof for the sub-Gaussian case, we use Vershynin's sub-Gaussian result for obtaining all the concentration bounds. However we can try to replace this by an approach motivated by the proof of [16, Theorem 4.1] that relies on the intuition that sub-Gaussian r.v.'s are bound whp. Thus, one can first work with a truncated sub-Gaussian, in our case, truncate each entry to $\sqrt{\log n}$ and use matrix Bernstein, and then deal with the extra errors introduced by this truncation. With this approach, we will get an extra factor of $n^{-c}$ in the SE bound. For large enough $n$, this extra factor is negligible. The advantage will be that we may only need $\alpha \geq Cr_v\log^3 n$ instead of $\alpha \geq Cn$.*

## 2.3 Result for only data-dependent noise case

**Corollary 2.12** (Only data-dependent noise)**.** *Given data vectors $\boldsymbol{y}_t := \boldsymbol{\ell}_t + \boldsymbol{w}_t$ with $\boldsymbol{w}_t = \boldsymbol{M}_t\boldsymbol{\ell}_t$, $t = 1, 2, \ldots, \alpha$. Let $\hat{\boldsymbol{P}}$ denote the matrix of top $r$ eigenvectors of $\boldsymbol{D} := \frac{1}{\alpha}\sum_t \boldsymbol{y}_t\boldsymbol{y}_t'$ and, for a scalar $q$, let*

$$d(\alpha) := cqf\sqrt{\frac{r\log n}{\alpha}} \text{ and } d_{\text{denom}}(\alpha) := cf\sqrt{\frac{r+\log n}{\alpha}}.$$

1. *If Assumption 2.4 holds, $\alpha^3 > (r \log n)$, and if, for scalars $q, b < 1$ satisfying $3\sqrt{b}qf + d(\alpha) + d_{\mathrm{denom}}(\alpha) < 1$, the matrices $\boldsymbol{M}_t$ can be decomposed as $\boldsymbol{M}_t = \boldsymbol{M}_{2,t}\boldsymbol{M}_{1,t}$ with $\boldsymbol{M}_{1,t}$ being such that*

$$\max_t \|\boldsymbol{M}_{1,t}\boldsymbol{P}\|_2 \le q \tag{4}$$

*and $\boldsymbol{M}_{2,t}$ being such that $\|\boldsymbol{M}_{2,t}\|_2 \le 1$ but*

$$\left\|\frac{1}{\alpha}\sum_{t=1}^{\alpha} \boldsymbol{M}_{2,t}\boldsymbol{M}_{2,t}'\right\|_2 \le b < 1, \tag{5}$$

*then, w.p. at least $1 - 10n^{-10}$,*

$$\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \le \frac{\sqrt{b}(2q + q^2)f + d(\alpha)}{1 - \sqrt{b}(2q + q^2)f - d(\alpha) - d_{\mathrm{denom}}(\alpha)}.$$

2. *If Assumption 2.5 holds (instead of Assumption 2.4), then we have the same result as above but with $d(\alpha) = d(\alpha)_{sG} := cf\sqrt{\frac{n}{\alpha}}$.*

*Proof Outline.* This is a corollary of Theorem 2.13. The proof idea is similar to that of Corollary 2.7. In this case we apply the $\sin\theta$ theorem with $\boldsymbol{D}_0 = \boldsymbol{P}(\frac{1}{\alpha}\sum_t \boldsymbol{a}_t\boldsymbol{a}_t')\boldsymbol{P}'$. Also, we need to carefully bound $\|\frac{1}{\alpha}\sum_{t=1}^{\alpha} \mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{w}_t']\|_2$ and $\|\frac{1}{\alpha}\sum_{t=1}^{\alpha} \mathbb{E}[\boldsymbol{w}_t\boldsymbol{w}_t']\|_2$; this is done in (7) below. $\boxtimes$

*Data-dependent noise.* Observe a few things about data-dependent noise. First, it is clearly non-isotropic and hence a condition that ensures that the noise is small compared to $\lambda^-$ is needed. Second, because of the assumed linear dependency on $\boldsymbol{\ell}_t$, and because (4) holds, the noise power depends linearly on maximum signal power, $\lambda^+$: we have $\|\mathbb{E}[\boldsymbol{w}_t\boldsymbol{w}_t']\|_2 \le q^2\lambda^+$. Here $q^2$ can be interpreted as a bound on the *noise-to-signal ratio*. Third, the signal-noise correlation is nonzero. In fact, its bound also linearly depends on $\lambda^+$: we have $\|\mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{w}_t']\|_2 \le q\lambda^+$. Since $q < 1$, the latter is, in fact, larger than the noise power bound. From the proof outline, for large enough $\alpha$, the subspace error bound essentially depends on the ratio $\|\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2/\lambda^-$. Since the signal-noise correlation is nonzero, this ratio is now bounded by $2\|\frac{1}{\alpha}\sum_{t=1}^{\alpha} \mathbb{E}[\boldsymbol{\ell}_t\boldsymbol{w}_t']\|_2 + \|\frac{1}{\alpha}\sum_{t=1}^{\alpha} \mathbb{E}[\boldsymbol{w}_t\boldsymbol{w}_t']\|_2$ instead of just the second term in the only uncorrelated noise case. Thus, without an assumption such as (5) on the signal-noise correlation, one would require $(2q + q^2)\lambda^+$ to be smaller than $0.45\epsilon\lambda^-$ to achieve subspace error below $\epsilon$. This is a hard requirement since it implies that $\epsilon$ can never be made smaller than the noise level, $q$.

Assuming (5) resolves the above issue. Observe that the subspace error depends on the time-averaged signal-noise correlation and time-averaged noise power, and not on their instantaneous values. The assumption (5) ensures that the bounds on the time-averaged values of both these are at least $\sqrt{b}$ times smaller than the bounds on their instantaneous values: by a careful application

of Cauchy-Schwartz inequality, it is not hard to see that (see (11) and (12) in the Appendix):

$$\left\| \frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}[\boldsymbol{\ell}_t \boldsymbol{w}_t{}'] \right\|_2 \le \sqrt{b} q \lambda^+, \text{ and} \tag{6}$$

$$\left\| \frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}[\boldsymbol{w}_t \boldsymbol{w}_t{}'] \right\|_2 \le \sqrt{b} q^2 \lambda^+. \tag{7}$$

Thus, because (5) holds, to achieve error $\epsilon$, other than a sample complexity lower bound, our result just needs $\sqrt{b}(2q + q^2)\lambda^+ < 0.45\epsilon\lambda^-$. With this, it is possible to achieve subspace recovery error $\epsilon$ that is smaller than $q$ by assuming that $b$ is small enough. Of course a lower bound on $\alpha$ that ensures $d(\alpha) < 0.1\epsilon$ and $cf\sqrt{\frac{r+\log n}{\alpha}} < 0.01$ will also be needed (discussed below).

*Examples where* (4) *and* (5) *holds.* One class of example situations where (5) would hold is when the data-dependent noise $\boldsymbol{w}_t$ is sparse (PCA in sparse data-dependent noise). Let $\mathcal{T}_t$ denote its support set. Then, in this case, $\boldsymbol{w}_t = \boldsymbol{I}_{\mathcal{T}_t} \boldsymbol{M}_{s,t} \boldsymbol{\ell}_t$ where $\boldsymbol{M}_{s,t}$ is a $|\mathcal{T}_t| \times n$ data-dependency matrix. If we pick $\boldsymbol{M}_{2,t} = \boldsymbol{I}_{\mathcal{T}_t}$, then $\sum_t \boldsymbol{M}_{2,t} \boldsymbol{M}_{2,t}{}'$ will be a diagonal matrix with $(i,i)$-th entry being equal to the number of time instants $t$ for which the index $i$ is part of the support $\mathcal{T}_t$. Hence $b$ will equal the maximum fraction of non-zeros in any row of the matrix $[\boldsymbol{w}_1, \boldsymbol{w}_2, \dots, \boldsymbol{w}_\alpha]$. Thus, in this case, (5) holds as long as this fraction is smaller than one. It holds with a small enough $b$ if this fraction is small enough. Moreover, (4) will hold as long as $\|\boldsymbol{M}_{1,t}\boldsymbol{P}\|_2 = \|\boldsymbol{M}_{s,t}\boldsymbol{P}\|_2$ is small. Two examples where this happens are given next.

A special case of the above problem is PCA in missing data. Let $\mathcal{T}_t$ denote the set of missing entries at time $t$. By setting the missing entries to zero, we can write out the observed data vector as $\boldsymbol{y}_t = \boldsymbol{\ell}_t - \boldsymbol{I}_{\mathcal{T}_t} \boldsymbol{I}_{\mathcal{T}_t}{}' \boldsymbol{\ell}_t$. Thus, in this case, $\boldsymbol{M}_{s,t} = -\boldsymbol{I}_{\mathcal{T}_t}{}'$ and so, $q$ is a bound on $\|\boldsymbol{I}_{\mathcal{T}_t}{}'\boldsymbol{P}\|_2$. Thus, for PCA-missing, $q$ will be small if columns of $\boldsymbol{P}$ are dense vectors and the number of missing entries at each time, $|\mathcal{T}_t|$, is small. We discuss this case further in Sec. 4.1. Another special case occurs in the subspace update step of ReProCS for dynamic robust PCA [9]. We explain this in detail in Sec. 4.2. Briefly, in this case, $\boldsymbol{M}_{s,t} = \boldsymbol{B}(\boldsymbol{I} - \hat{\boldsymbol{P}}\hat{\boldsymbol{P}}')$ where $\boldsymbol{B}$ is a matrix satisfying $\|\boldsymbol{B}\| \le 1.2$ and $\hat{\boldsymbol{P}}$ is a previous estimate of span$(\boldsymbol{P})$ satisfying SE$(\hat{\boldsymbol{P}}, \boldsymbol{P}) < q/1.2 \ll 1$.

Another related class of problems where the above assumptions would hold is if $\boldsymbol{w}_t$ is sparse in a basis or dictionary $\boldsymbol{Q}$. Then, $\boldsymbol{w}_t = \boldsymbol{Q}\boldsymbol{I}_{\mathcal{T}_t} \boldsymbol{M}_{s,t} \boldsymbol{\ell}_t$. In this we can use $\boldsymbol{M}_{2,t} = \boldsymbol{Q}\boldsymbol{I}_{\mathcal{T}_t}/\|\boldsymbol{Q}\|_2$ and $\boldsymbol{M}_{1,t} = \|\boldsymbol{Q}\|_2\boldsymbol{M}_{s,t}$.

*Sample complexity.* To achieve error below $\epsilon$, in the bounded case, the sample complexity needed is $\alpha \ge C \max(\frac{q^2 f^2}{\epsilon^2}(r \log n), f^2(r+\log n))$. Thus, if $\epsilon$ is a constant fraction of $q$, the sample complexity needed is just $\alpha \ge Cf^2(r \log n)$. This is nearly optimal. For constant $f$, this is only $O(\log n)$ times the minimum number of samples required to even define an $r$-dimensional subspace. In the sub-Gaussian case, according to our current result, $O(n)$ samples are needed, however as explained earlier, this can possibly be improved in certain settings.

*Comparing Corollaries 2.7 and 2.12.* In both results, the bound on SE depends on the condition number $f = \lambda^+/\lambda^-$, however, the dependence is much weaker in the uncorrelated noise case. In this case, $f$ only appears in terms that contain the sample complexity $\alpha$. Thus, any $f$ can be dealt with by picking a proportionally larger $\alpha$. However in the data-dependent noise case, $f$ also appears in a term other than $d(\alpha)$ or $d_{\text{denom}}(\alpha)$. To get SE below $\epsilon$, in this case, one needs $\sqrt{b}(2q+q^2)f < 0.45\epsilon$. This is a much stronger requirement since it cannot be ensured just by using more samples $\alpha$. This is needed because the data-dependent noise power depends on $\lambda^+$.

On the other hand, consider the sample complexity $\alpha$ required to achieve error $\epsilon$ that is a constant fraction of the noise level in the bounded data and noise setting. In the uncorrelated noise case, this is $C \max(r_v, r) \log n$ whereas in the data-dependent noise case, this is only $C(r \log n)$. If $r_v$ is larger than $r$, the former will need more samples.

## 2.4 The general result

We now give the most general result.

**Theorem 2.13.** *Given data vectors $\boldsymbol{y}_t := \boldsymbol{\ell}_t + \boldsymbol{w}_t + \boldsymbol{v}_t$ with $\boldsymbol{w}_t = \boldsymbol{M}_t \boldsymbol{\ell}_t$ and $\boldsymbol{v}_t$ and $\boldsymbol{\ell}_t$ uncorrelated. Let $\hat{\boldsymbol{P}}$ denote the matrix of top $r$ eigenvectors of $\boldsymbol{D} := \frac{1}{\alpha} \sum_t \boldsymbol{y}_t \boldsymbol{y}_t'$ and define*

$$d(\alpha) := c\sqrt{\eta} \max \left( qf\sqrt{\frac{r\log n}{\alpha}}, g\sqrt{\frac{\max(r_v, r)\log n}{\alpha}} \right),$$

*and*

$$d_{\text{denom}}(\alpha) := c\eta f \sqrt{\frac{r + \log n}{\alpha}}.$$

1. *If Assumption 2.4 holds, $\alpha^3 > \max(r_v, r) \log n$,*

   (a) *if, for a $b < 1$ and a $q < 1$, the data-dependency matrices $\boldsymbol{M}_t$ satisfy the assumption given in Corollary 2.12,*

   (b) *and if $\frac{\lambda^+_{v,\text{rest}} - \lambda^-_{v,\boldsymbol{P}}}{\lambda^-} + 3\sqrt{b}qf + d(\alpha) + d_{\text{denom}}(\alpha) < 1$,*

   *then, w.p. at least $1 - 10n^{-10}$,*

   $$\text{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq \frac{\frac{\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}}{\lambda^-} + \sqrt{b}(2q + q^2)f + d(\alpha)}{1 - \frac{\lambda^+_{v,\text{rest}} - \lambda^-_{v,\boldsymbol{P}}}{\lambda^-} - \sqrt{b}(2q + q^2)f - d(\alpha) - d_{\text{denom}}(\alpha)}.$$

2. *If Assumption 2.5 holds (instead of Assumption 2.4), then we have the same result as above but with $d(\alpha) = d(\alpha)_{sG} := c \max \left( \frac{\lambda^+_v}{\lambda^-}, f \right) \sqrt{\frac{n}{\alpha}}$ and we do not need $\alpha^3 \geq r \log n$. Also, the result now holds w.p. greater than $1 - 10 \exp(-cn)$.*

11

*Proof.* See Appendices A and C for the bounded and sub-Gaussian cases respectively.

Corollary 2.7 follows from the above result by setting $q = 0$. Corollary 2.12 follows by setting $r_v = 0$ and $\Sigma_v^+ = 0$. In both corollaries, we treat $\eta$ as a numerical constant.

From Theorem 2.13, to get error below $\epsilon$ in the most general case, we need a sample complexity lower bound *and* we the signal-noise correlation assumption stated in Corollary 2.12 to hold with parameters $b, q$ that satisfy

$$\sqrt{b}(2q + q^2)f + \frac{\lambda_{v,\mathrm{rest}}^+ - \lambda_{v,\boldsymbol{P}}^-}{\lambda^-} < 0.5, \text{ and } \frac{\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}}{\lambda^-} + \sqrt{b}(2q + q^2)f < 0.25\epsilon.$$

In the bounded case, the sample complexity required is
$\alpha \geq C \max\left(\frac{q^2}{\epsilon^2} \max(r_v, r) \log n, \frac{(qf)^2}{\epsilon^2} r \log n, f^2(r + \log n)\right)$. In the sub-Gaussian case, we need $\alpha \geq C \frac{q^2}{\epsilon^2} n$ although, as discussed earlier in Remark 2.11, this can possibly be improved.

**Remark 2.14.** *If $\boldsymbol{\Lambda}$ and $\boldsymbol{\Sigma}_v$ were time-varying, the above result will hold with the following simple changes.*

*(1) Define "average" versions of $\lambda^-$ and $\lambda_{v,\boldsymbol{P}}^-$ as $\bar{\lambda}^- := \lambda_{\min}(\frac{1}{\alpha} \sum_t \boldsymbol{\Lambda}_t)$. Define $\bar{\lambda}_{v,\boldsymbol{P}}^- := \lambda_{\min}(\boldsymbol{P}'(\frac{1}{\alpha} \sum_t \boldsymbol{\Sigma}_{v,t})\boldsymbol{P})$. In the result above, replace $\lambda^-$ and $\lambda_{v,\boldsymbol{P}}^-$ by their "average" versions $\bar{\lambda}^-$ and $\bar{\lambda}_{v,\boldsymbol{P}}^-$ respectively.*

*(2) Define $\lambda_{max}^+ := \max_t \lambda_{\max}(\boldsymbol{\Lambda}_t)$. Similarly define "max" versions of $\lambda_v^+$, $\lambda_{v,\mathrm{rest}}^+$, and $\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}$. In the result above, replace $\lambda^+$, $\lambda_v^+$, $\lambda_{v,\mathrm{rest}}^+$, $\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}$ by their "max" versions $\lambda_{max}^+$, $\lambda_{v,max}^+$, $\lambda_{v,\mathrm{rest},max}^+$, $\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp,max}$ respectively.*

*(3) Because of the above two changes, $f$ gets replaced by $\lambda_{max}^+/\bar{\lambda}^-$.*

## 2.5 Automatically estimating $r$

In the above result and its corollaries, we assumed that $r$, which is the signal subspace dimension, is known. In practice however this is usually unknown. There are two easy and commonly used ways to automatically estimate $r$. The first approach is as done in [6]. This computes $\hat{r}$ as the smallest index $j$ for which the $j$-th eigenvalue of $\boldsymbol{D} := \sum_{t=1}^\alpha \boldsymbol{y}_t \boldsymbol{y}_t'$ is above a threshold. Thus,

$$\hat{r} := \arg\min\{j : \lambda_j(\boldsymbol{D}) \geq 0.5\lambda^-\}. \tag{8}$$

Notice that this requires knowledge of $\lambda^-$. However, as we will see, this does not require extra assumptions beyond what Theorem 2.13 already assumes. An alternate way to estimate $r$ is by looking for the largest eigen-gap, i.e.,

$$\hat{r} := \arg\max_j[\lambda_j(\boldsymbol{D}) - \lambda_{j+1}(\boldsymbol{D})]. \tag{9}$$

This does not require knowledge of any model parameter. However, as we see below, this works only under the assumption that consecutive eigenvalues of the

matrix $\boldsymbol{\Lambda} + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}$ do not have a large gap. It is also more expensive since it requires computing all eigenvalues of $\boldsymbol{D}$.

Consider (8). To prove that this works, we need to show that $\lambda_r(\boldsymbol{D}) \geq 0.5\lambda^-$ and $\lambda_{r+1}(\boldsymbol{D}) < 0.5\lambda^-$. Let $\boldsymbol{D}_0 = \boldsymbol{P}(\boldsymbol{\Lambda} + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P})\boldsymbol{P}'$. By Weyl, $\lambda_r(\boldsymbol{D}) \geq \lambda_r(\boldsymbol{D}_0) - \|\boldsymbol{D} - \boldsymbol{D}_0\|_2 \geq \lambda^- + \lambda^-_{v,\boldsymbol{P}} - \|\boldsymbol{D} - \boldsymbol{D}_0\|_2$ and $\lambda_{r+1}(\boldsymbol{D}) \leq 0 + \|\boldsymbol{D} - \boldsymbol{D}_0\|_2$. Using (10), (11), (12), and Lemma A.20 from the Appendix,

$$\|\boldsymbol{D} - \boldsymbol{D}_0\|_2 \leq \Delta\lambda^-, \text{ where } \Delta := d_{\mathrm{denom}}(\alpha) + d(\alpha) + 3\sqrt{b}qf + \frac{\lambda^+_{v,\mathrm{rest}}}{\lambda^-}$$

where $d_{\mathrm{denom}}(\alpha), d(\alpha)$ are defined in Theorem 2.13. Thus, we have the following result.

**Theorem 2.15** (Estimating $r$ using (8))**.** *Assume that Assumption 2.4 holds and the assumption on $\boldsymbol{M}_t$ given in Corollary 2.12 holds. Let $\Delta$ be as defined above. If $\Delta < \frac{1}{2}$, then w.p. at least $1 - 10n^{-10}$, (8) returns the correct estimate of $r$. This result also holds if Assumption 2.4 is replaced by Assumption 2.5 as long as we replace $d(\alpha)$ by $d(\alpha)_{sG}$.*

Proceeding as above for (9),

$\lambda_r(\boldsymbol{D}) - \lambda_{r+1}(\boldsymbol{D}) \geq \lambda^- + \lambda^-_{v,\boldsymbol{P}} - 2\Delta\lambda^-$,

for $j < r$, $\lambda_j(\boldsymbol{D}) - \lambda_{j+1}(\boldsymbol{D}) \leq (\lambda_j(\boldsymbol{\Lambda} + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}) - \lambda_{j+1}(\boldsymbol{\Lambda} + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P})) + 2\Delta\lambda^-$,

for $j > r$, $\lambda_j(\boldsymbol{D}) - \lambda_{j+1}(\boldsymbol{D}) \leq 2\Delta\lambda^-$

Thus we the following result for (9).

**Theorem 2.16** (Estimating $r$ using (9))**.** *Assume that Assumption 2.4 holds and the assumption on $\boldsymbol{M}_t$ given in Corollary 2.12 holds. Let $\Delta$ be as defined above. If*

$$\max_{j<r}(\lambda_j(\boldsymbol{\Lambda} + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}) - \lambda_{j+1}(\boldsymbol{\Lambda} + \boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P})) \leq (1 - 4\Delta)\lambda^- + \lambda^-_{v,\boldsymbol{P}}$$

*then w.p. at least $1 - 10n^{-10}$, (9) returns the correct estimate of $r$. This result also holds if Assumption 2.4 is replaced by Assumption 2.5 as long as we replace $d(\alpha)$ by $d(\alpha)_{sG}$.*

# 3 Discussion of Related Work

A detailed discussion is given here.

*Discussion of [2].* This work was the first to obtain finite sample guarantees for PCA. Its main result, [2, Theorem 2.1], assumes a spiked covariance model with $r = 1$ spike and Gaussianity of both data and noise. It was proved using a different set of concentration bounds and hence its exact form is a little different from our result in this setting. However, if one looks at the dominant terms in its required assumption or in its upper bound on $\sin\theta_{PCA}$, the conclusions are

the same as those of our Corollary 2.9. In our notation, $\sin\theta_{PCA} \equiv \mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P})$. First to explain notation, its $p \equiv n$, $n \equiv \alpha$, $\kappa^2 \approx \lambda^-$ (actually $\mathbb{E}[\kappa^2] \equiv \lambda^-$), $\sigma_v^2 \equiv \lambda_v^+$. In our notation, [2, Theorem 2.1] says the following. When $n \geq \alpha$, if $\lambda^- \gtrsim \lambda_v^+ \frac{n}{\alpha}$, then, whp, $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \lesssim c\sqrt{\frac{\lambda_v^+}{\lambda^-}}\sqrt{\frac{n}{\alpha}}$. Here $\gtrsim, \lesssim$ indicate that we are only using the dominant terms from their long expression.

Consider our Corollary 2.9 with $r = 1$, Gaussian data and noise, and $n \geq \alpha$. Since $r = 1$, so $\lambda^+ = \lambda^-$, $f = 1$, and $g = \max(\frac{\lambda_v^+}{\lambda^-}, \sqrt{\frac{\lambda_v^+}{\lambda^-}})$. Ignoring constants, Corollary 2.9 assumes $g < \sqrt{\frac{\alpha}{n}}$. Since $\sqrt{\frac{\alpha}{n}} < 1$, the max in the $g$ expression is achieved by the square root term. Thus, in this setting, Corollary 2.9 says the following: if $2\sqrt{\frac{\lambda_v^+}{\lambda^-}}\sqrt{\frac{n}{\alpha}} < 1$ (equivalently $\lambda^- \geq 4\lambda_v^+ \frac{n}{\alpha}$), then, w.p. at least $1 - 10\exp(-cn)$, $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq 2\sqrt{\frac{\lambda_v^+}{\lambda^-}}\sqrt{\frac{n}{\alpha}}$. This is the same as the simplified version of [2, Theorem 2.1] given above.

Because [2] only considered the $r = 1$ and spiked covariance model setting, it was able to provide more insight into its guarantees beyond just an SE upper bound. It showed that its upper bound on subspace error is sharp by also providing an expression for the expected subspace error. Moreover, it provided an approximate expression for the top eigenvector of the sample covariance matrix that is valid when the noise variance is small. For $r > 1$, these things are difficult to do. For the setting in our paper (non-isotropic and possibly data-dependent noise), these are even harder to do. We only give an example, Example 2.10, to show that, the subspace error will not be small if a bound on noise power outside span($\boldsymbol{P}$) is not assumed.

*Comparing Theorem 2.13 with the result of [6].* Our result is a significant improvement over that of [6] where the correlated-PCA problem was first studied. We include a second uncorrelated noise component in our result which makes the data model more practically valid. Second, we also get results under a general sub-Gaussian data and noise assumption.

To compare with the result of [6], consider Corollary 2.12 under the bounded assumption. The signal-noise correlation model assumed in it is a significant simplification of the one needed by the result of [6]. That result needed $\|\frac{1}{\alpha}\sum_{t=1}^\alpha \boldsymbol{M}_{2,t}\boldsymbol{A}_t\boldsymbol{M}_{2,t}'\|_2 \leq b$ to hold for all sets of positive semi-definite (p.s.d.) matrices $\boldsymbol{A}_t$, $t = 1, 2, \ldots, \alpha$. This is a much stronger requirement. Our current result only needs this to hold only for $\boldsymbol{A}_t = \boldsymbol{I}$, i.e., it needs (5) to hold. Consider the sparse $\boldsymbol{w}_t$ example. For this, as explained earlier, (5) would hold if the fraction of nonzeros in any row of the noise matrix $[\boldsymbol{w}_1, \boldsymbol{w}_2, \ldots, \boldsymbol{w}_\alpha]$ is bounded by $b$. On the other hand, it is not clear if the assumption needed by [6] holds for this example. The examples given in [6] involved much more stringent assumptions on $\mathcal{T}_t$ - the sets $\mathcal{T}_t$ needed to be either mutually disjoint, or mutually disjoint every few frames, or they needed to change in a way to model stop and go object motion in one direction.

Second, our sample complexity bound is a significant improvement over that of [6]. If the desired error $\epsilon$ is much larger than $qf$, our required sample complexity is $O(r + \log n)$. If $r \geq c\log n$, this is optimal. In the more common small $\epsilon$ setting, to get the subspace error to below say $\epsilon = q/4$, we

need $\alpha \geq 16Cf^2(r \log n)$ samples. This is also much better than the earlier bound from [6] of $\alpha \geq C\frac{f^2}{\epsilon^2}(r^2 \log n)$ which implied that $\alpha \geq 16C\frac{f^2}{q^2}(r^2 \log n)$ was needed to achieve the above subspace error level. This older bound had an extra multiplicative factor of $r/q^2$. In our work, we remove the extra $r$ factor by using matrix Bernstein to replace matrix Hoeffding to get high probability bounds on the deviation between time-averaged signal-noise correlation and noise power and their respective expected values. We remove the extra $1/q^2$ factor by bounding the $r$-th eigenvalue of $\sum_t \boldsymbol{\ell}_t \boldsymbol{\ell}_t{}' = \boldsymbol{P}(\sum_t \boldsymbol{a}_t \boldsymbol{a}_t{}')\boldsymbol{P}'$ by using the sub-Gaussian result of Vershynin (Theorem 5.39 of [13]) to bound the minimum eigenvalue of $\sum_t \boldsymbol{a}_t \boldsymbol{a}_t{}'$. In [6], the authors had used matrix Hoeffding for this term as well.

*Discussion of [11].* In [11] and references therein, the authors study the effect of multiplicative perturbations of Hermitian matrices on their principal subspaces. This line of work provides a tighter bound than Davis-Kahan for the subspace error between principal subspaces of a Hermitian matrix $A$ and of its perturbed version $BAB'$ for a non-singular matrix $B$. However, such results are not applicable for our problem even in the only data-dependent noise case, since $\boldsymbol{w}_t$ satisfies $\boldsymbol{w}_t = \boldsymbol{M}_t \boldsymbol{\ell}_t$ where $\boldsymbol{M}_t$ is time-varying.

*Discussion of [12].* This work develops concentration inequalities for sample covariance matrices in a setting where observed data lie in a Euclidean ball of radius $O(\sqrt{n})$. It obtains new results for the setting where the observed data is bounded but is not a "nice" sub-Gaussian. They explain that, if observed data satisfies their equation (1.5), then $O(n \log n)$ samples are needed to ensure that the sample covariance matrix is close to its expected value. This matches the sample complexity predicted by our result for the bounded case with $r_v = n$ and without any other assumption. In this case our observed data $\boldsymbol{y}_t = \boldsymbol{\ell}_t + \boldsymbol{w}_t + \boldsymbol{v}_t$ satisfies equation (1.5) of that paper with $K^2 = \eta\lambda^+(1+q^2)+\lambda_v^+$ and $L^2 = cK^2$. When $r_v = Cn$, if we add a mild extra assumption that the $k$-th moment of $\boldsymbol{v}_t$ for a $k > 4$ is bounded, then we can tap into the main result of [12] to show that the sample complexity can be reduced to from $O(n \log n)$ to $O(n(\log \log n)^2)$.

# 4   Application to PCA in sparse data-dependent noise and its special cases

Consider the PCA in sparse data-dependent noise (PCA-SDDN) problem described earlier. In this case, $\boldsymbol{y}_t = \boldsymbol{\ell}_t + \boldsymbol{w}_t$ with $\boldsymbol{w}_t = \boldsymbol{I}_{\mathcal{T}_t}\boldsymbol{M}_{s,t}\boldsymbol{\ell}_t$. Thus $\boldsymbol{w}_t$ is sparse with support $\mathcal{T}_t$. The following is an easy corollary of Corollary 2.12.

**Corollary 4.17** (PCA-SDDN). *Given data vectors $\boldsymbol{y}_t := \boldsymbol{\ell}_t + \boldsymbol{I}_{\mathcal{T}_t}\boldsymbol{M}_{s,t}\boldsymbol{\ell}_t$, $t = 1, 2, \ldots, \alpha$. Let $\hat{\boldsymbol{P}}$ be the matrix of top $r$ eigenvectors of $\boldsymbol{D} := \frac{1}{\alpha}\sum_t \boldsymbol{y}_t\boldsymbol{y}_t'$. Assume that $\boldsymbol{\ell}_t$ satisfies Assumption 2.4, $\max_t \|\boldsymbol{M}_{s,t}\boldsymbol{P}\|_2 \leq q < 1$, the fraction of nonzeroes in any row of the noise matrix $[\boldsymbol{w}_1, \boldsymbol{w}_2, \ldots, \boldsymbol{w}_\alpha]$ is bounded by $b$, and $b, q$ satisfy $3\sqrt{b}qf + d(\alpha) + d_{\mathrm{denom}}(\alpha) < 1$. Here $d(\alpha), d_{\mathrm{denom}}(\alpha)$ are as defined in Corollary 2.12. Then, w.p. at least $1 - 10n^{-10}$, $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq$*

$$\frac{3\sqrt{b}qf+d(\alpha)}{1-(3\sqrt{b}qf+d(\alpha)+d_{\text{denom}}(\alpha))}.$$

**Remark 4.18** (PCA-SDDN - alternate). *Another way to state Corollary 4.17 is as follows. Assume that $\boldsymbol{\ell}_t$ satisfies Assumption 2.4, $\max_t \|\boldsymbol{M}_{s,t}\boldsymbol{P}\|_2 \leq q < 1$, and the fraction of nonzeroes in any row of the noise matrix $[\boldsymbol{w}_1, \boldsymbol{w}_2, \ldots, \boldsymbol{w}_\alpha]$ is bounded by b. For an $\epsilon_{\text{SE}} > 0$, if $\alpha \geq \alpha_0 = C \max \left( \frac{q^2 f^2}{\epsilon_{\text{SE}}^2}(r\log n), f^2(r+\log n) \right)$ and if $3\sqrt{b}qf < 0.9\frac{\epsilon_{\text{SE}}}{1+\epsilon_{\text{SE}}}$, then w.p. at least $1 - 10n^{-10}$, $\text{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq \epsilon_{\text{SE}}$.*

It is also possible to solve the above problem using techniques from the sparse + low-rank matrix recovery (robust PCA) literature, e.g., [17, 18, 19, 20]. These also do not assume anything about whether the noise (sparse outlier) depends on the true data or not and hence do allow data-dependent noise. However, there are some differences. (1) Our guarantee for PCA via SVD (and, in fact, all guarantees for PCA) *assume* that the noise is smaller than than the data (needs $q < 1$) where as robust PCA solutions are designed to handle noise (sparse outliers) that can have any magnitude. (2) Because of this, the robust PCA solutions are more expensive than the simple SVD solution that works for PCA. The most recent robust PCA solutions [19, 20] have nearly the same order of complexity as simple SVD, however in practice they are still slower. (3) More importantly, all robust PCA via sparse + low-rank recovery solutions require the columns of $\boldsymbol{P}$ to be dense (not sparse). Their guarantees also require denseness of the right singular vectors. In our notation, these would be columns of the matrix $[\boldsymbol{a}_1, \boldsymbol{a}_2, \ldots, \boldsymbol{a}_\alpha]'\boldsymbol{\Lambda}^{-1}$. These assumptions are necessary even for identifiability; otherwise the true data vector(s) may get wrongly classified as sparse outliers by the algorithm. If the goal is only PCA (and not recovering the low rank matrix), then it is possible that the denseness assumption on right singular vectors can be removed, see e.g. [21]. Our guarantee given above for simple SVD does not require denseness of even the columns of $\boldsymbol{P}$. As shown in [6, Table 1], when $\boldsymbol{P}$ contains sparse vectors, robust PCA solutions also fail in practice, while SVD does not (as long as the sparse noise is small of course).

The PCA-SDDN model given above is often a valid one for video analytics applications, where $\boldsymbol{\ell}_t$ is the background layer of image frame $t$, $\mathcal{T}_t$ is the foreground (e.g., occlusion) support of frame $t$, and $\boldsymbol{w}_t := \boldsymbol{I}_{\mathcal{T}_t}\boldsymbol{M}_{s,t}\boldsymbol{\ell}_t$ is the difference between foreground and background intensities on $\mathcal{T}_t$. The above corollary is useful in problems involving subspace learning of slow changing videos (well modeled as being low rank) when the video is corrupted by foreground occlusions whose intensity is very similar to that of the background and that are correlated with the background (resulting in small magnitude and data-dependent $\boldsymbol{w}_t$). Occlusions due to shadows often fall in this category. Another application is in using functional MRI (fMRI) data to learn the low-dimensional subspace in which resting state fMRI data lies. Even in the absence of external stimuli, it is well understood that the human brain is never fully resting. The sparse and small magnitude activations generated by random thoughts in the so-called "resting state" brain are well modeled as $\boldsymbol{w}_t$ described above. Besides these examples, we discuss two other practically relevant special cases of PCA-SDDN in the next two subsections.

## 4.1 Special case: PCA with missing data

Let $\mathcal{T}_t$ denote the set of missing entries at time $t$. As explained earlier, if we set the missing entries to zero to define an $n$-length observed data vector $\boldsymbol{y}_t$, then $\boldsymbol{y}_t$ satisfies $\boldsymbol{y}_t = \boldsymbol{\ell}_t + \boldsymbol{w}_t$ where $\boldsymbol{w}_t = -\boldsymbol{I}_{\mathcal{T}_t}\boldsymbol{I}_{\mathcal{T}_t}'\boldsymbol{\ell}_t$ is the sparse "error" or "noise" due to the missing part of the data. In this case, $\boldsymbol{M}_{s,t} = -\boldsymbol{I}_{\mathcal{T}_t}'$. The bound on $q$ thus translates to a denseness assumption on the columns of $\boldsymbol{P}$. Let $\mu$ be the densenesss (incoherence) parameter [17] for $\boldsymbol{P}$, i.e., let $\mu$ be the smallest real number so that

$$\max_i \|\boldsymbol{I}_i'\boldsymbol{P}\|_2^2 \le \mu^2 r/n.$$

Also let $s := \max_t |\mathcal{T}_t|$ be an upper bound on the number of missing entries at any time. It is easy to see that $\max_t \|\boldsymbol{M}_{s,t}\boldsymbol{P}\|_2^2 = \max_t \|\boldsymbol{I}_{\mathcal{T}_t}'\boldsymbol{P}\|_2^2 \le s\max_{i=1,2,\dots,n} \|\boldsymbol{I}_i'\boldsymbol{P}\|_2^2 \le \mu^2 rs/n \equiv q^2$. We have the following corollary.

**Corollary 4.19** (PCA with missing data). *Given data vectors $\boldsymbol{y}_t := \boldsymbol{\ell}_t - \boldsymbol{I}_{\mathcal{T}_t}\boldsymbol{I}_{\mathcal{T}_t}'\boldsymbol{\ell}_t$, $t = 1,2,\dots,\alpha$ with $|\mathcal{T}_t| \le s$. Assume that $\boldsymbol{\ell}_t$ satisfies Assumption 2.4, the fraction of missing entries in any row of the data matrix is at most $b$ where $b$ satisfies $3\sqrt{b}qf + d(\alpha) + d_{\mathrm{denom}}(\alpha) < 1$ and $q = \sqrt{\mu^2 rs/n}$. Here $d(\alpha), d_{\mathrm{denom}}(\alpha)$ are defined in Corollary 2.12. Then, w.p. at least $1 - 10n^{-10}$, $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \le \frac{3\sqrt{b}qf + d(\alpha)}{1 - 3\sqrt{b}qf - d(\alpha) - d_{\mathrm{denom}}(\alpha)}$.*

Another way to solve the above problem would be to use techniques from the low-rank matrix completion (LRMC) literature to first complete the incomplete low-rank true data matrix $\boldsymbol{L}$ and then compute its left singular vectors. (1) This will need fewer observed entries because it relies on more expensive techniques that are designed to actually deal with missing data instead of just SVD which treats the missing data as noise. However, the LRMC methods will also be much slower. In fact, SVD is often the initialization step for iterative LRMC solutions, e.g., [22]. (2) It is hard to directly compare Corollary 4.19 with the LRMC guarantees since both use different assumptions, but the following can be said. LRMC results assume that observed entries are selected uniformly at random (or via a Bernoulli model), while Corollary 4.19 assumes a bound on the number of missing entries per row ($b$) and per column ($s/n$). (3) Moreover, LRMC results require denseness of left and right singular vectors, while our Corollary 4.19 only needs denseness of left singular vectors.

## 4.2 Special case: subspace update step of a dynamic robust PCA solution

A very important special case of PCA-SDDN occurs in the subspace update step of the Recursive Projected Compressive Sensing (ReProCS) approach to dynamic robust PCA [8, 9]. For dynamic robust PCA, the observed data vector $\boldsymbol{m}_t$ satisfies $\boldsymbol{m}_t := \boldsymbol{\ell}_t + \boldsymbol{x}_t$, where $\boldsymbol{x}_t$ is a sparse outlier with support denoted by $\mathcal{T}_t$ at time $t$, and $\boldsymbol{\ell}_t := \boldsymbol{P}_t\boldsymbol{a}_t$ is the true data vector that lies in a low dimensional subspace that is "slowly" changing; the subspace is fixed for a while and then changes by a little. To be precise we assume that $\boldsymbol{P}_t = \boldsymbol{P}_{t_j}$ for $t \in [t_j, t_{j+1})$ with

$SE(\boldsymbol{P}_{t_{j-1}}, \boldsymbol{P}_{t_j}) \le b_P \ll 1$ and with $t_{j+1} - t_j$ is lower bounded. One generative model for this is given in [23, Equation (2)]. The goal is to track this changing subspace over time. The columns of the matrices $\boldsymbol{P}_{t_j}$'s are assumed to be dense (not sparse). For simplicity we often use $\boldsymbol{P}_j := \boldsymbol{P}_{t_j}$.

ReProCS proceeds as follows. Given an accurate estimate of the previous subspace, denoted $\hat{\boldsymbol{P}}_{t-1}$, it first projects $\boldsymbol{m}_t$ orthogonal to $\hat{\boldsymbol{P}}_{t-1}$ to get $\tilde{\boldsymbol{m}}_t := (\boldsymbol{I} - \hat{\boldsymbol{P}}_{t-1}\hat{\boldsymbol{P}}_{t-1}')\boldsymbol{m}_t$. Because of the slow subspace change assumption, it can be argued that this nullifies most of $\boldsymbol{\ell}_t$ and gives projected measurements of $\boldsymbol{x}_t$. The problem of recovering $\boldsymbol{x}_t$ from $\tilde{\boldsymbol{m}}_t$ is now a standard Compressive Sensing (CS) problem [24] in small noise, $\boldsymbol{\beta}_t := (\boldsymbol{I} - \hat{\boldsymbol{P}}_{t-1}\hat{\boldsymbol{P}}_{t-1}')\boldsymbol{\ell}_t$. ReProCS uses ell-1 minimization followed by support estimation and Least Squares based debiasing to solve this CS problem. Once $\boldsymbol{x}_t$ is recovered, it recovers $\boldsymbol{\ell}_t$ by subtraction, $\hat{\boldsymbol{\ell}}_t = \boldsymbol{m}_t - \hat{\boldsymbol{x}}_t$. The estimates $\hat{\boldsymbol{\ell}}_t$ are used to update the subspace estimate every $\alpha$ frames by solving either a PCA or an incremental PCA problem. It is assumed that the subspace change is slow enough so that $t_{j+1} - t_j > K\alpha$. This allows the subspace to be updated $K$ times, each time with a new set of $\alpha$ frames of $\hat{\boldsymbol{\ell}}_t$, before it changes. The intuitive reason why this works is, after each subspace update, $\boldsymbol{\beta}_t$ reduces and hence the CS step error reduces. Thus, the error in $\hat{\boldsymbol{\ell}}_t$ reduces and this, in turn, helps reduce the subspace recovery error at the next update.

To understand in a simple fashion how PCA-SDDN fits in here, assume that the subspace change is detected exactly at $t_j$ and the subspace update times are also aligned so that the first subspace update is done at $t = t_j + \alpha - 1$, the second is at $t = t_j + 2\alpha - 1$, and so on. Let $\hat{\boldsymbol{P}}_{j,k}$ denote the updated subspace estimate after the $k$-th update with $\hat{\boldsymbol{P}}_{j,0} = \hat{\boldsymbol{P}}_{j-1}$. Thus, in the interval $[t_j, t_j + \alpha)$, $\boldsymbol{P}_t = \boldsymbol{P}_j$ and $\hat{\boldsymbol{P}}_{t-1} = \hat{\boldsymbol{P}}_{j,0} = \hat{\boldsymbol{P}}_{j-1}$. Assume that the previous subspace is recovered with $\epsilon$ error, i.e. $SE(\hat{\boldsymbol{P}}_{j-1}, \boldsymbol{P}_{j-1}) \le \epsilon$ before $t_j$. Using this and the denseness of columns of $\boldsymbol{P}_t$, one can show that $\boldsymbol{x}_t$ is recovered accurately for all $t \in [t_j, t_j + \alpha)$. The same analysis also shows that $\|\boldsymbol{B}_t\| \le 1.2$ where $\boldsymbol{B}_t := \boldsymbol{I}_{\mathcal{T}_t}'(\boldsymbol{I} - \hat{\boldsymbol{P}}_{j,0}\hat{\boldsymbol{P}}_{j,0}')\boldsymbol{I}_{\mathcal{T}_t}$. Then, using simple extra assumptions, one can argue that $\mathcal{T}_t$, which is the support of the outlier vector $\boldsymbol{x}_t$, can be recovered exactly. With this, it can be shown that $\boldsymbol{e}_t := \hat{\boldsymbol{\ell}}_t - \boldsymbol{\ell}_t = \boldsymbol{x}_t - \hat{\boldsymbol{x}}_t$ satisfies

$$\boldsymbol{e}_t = \boldsymbol{I}_{\mathcal{T}_t}\boldsymbol{B}_t^{-1}\boldsymbol{I}_{\mathcal{T}_t}'(\boldsymbol{I} - \hat{\boldsymbol{P}}_{j,0}\hat{\boldsymbol{P}}_{j,0}')\boldsymbol{\ell}_t.$$

Thus, the subspace update step is an instance of PCA-SDDN with $\boldsymbol{y}_t \equiv \hat{\boldsymbol{\ell}}_t$ and $\boldsymbol{w}_t \equiv \boldsymbol{e}_t$ for $t = t_j, t_j + 1, \ldots t_j + \alpha - 1$. We can apply Corollary 4.17 with $b$ being the maximum fraction of nonzeros in any row of $[\boldsymbol{x}_{t_j}, \boldsymbol{x}_{t_j+1}, \ldots, \boldsymbol{x}_{t_j+\alpha-1}]$ and with $\boldsymbol{M}_{s,t} = \boldsymbol{B}_t^{-1}\boldsymbol{I}_{T_t}'(\boldsymbol{I} - \hat{\boldsymbol{P}}_{j,0}\hat{\boldsymbol{P}}_{j,0}')$. Thus, $\|\boldsymbol{M}_{s,t}\boldsymbol{P}\|_2 \le 1.2SE(\hat{\boldsymbol{P}}_{j,0}, \boldsymbol{P}_j) = 1.2SE(\hat{\boldsymbol{P}}_{j-1}, \boldsymbol{P}_j) \le 1.2(\epsilon + b_P) := q_0$. Apply the PCA-SDDN result with $q = q_0$ and $\epsilon_{SE} = q_0/4$. Thus, if $3\sqrt{b}f < 0.2$, and if $\alpha \ge \alpha_0 = Cf^2(r \log n)$, then, after the first subspace update at $t = t_j + \alpha - 1$, $SE(\hat{\boldsymbol{P}}_{j,1}, \boldsymbol{P}_j) \le q_0/4$. This, in turn ensures that the bound on the noise $\boldsymbol{\beta}_t$ seen by the CS steps for the next interval, $[t_j + \alpha, t_j + 2\alpha - 1)$, is significantly smaller. Using denseness and the bound on $SE(\hat{\boldsymbol{P}}_{j,1}, \boldsymbol{P}_j)$ and simple extra assumptions, one can then argue

that the CS step gives a more accurate estimate of $\boldsymbol{x}_t$ and that $\mathcal{T}_t$ is correctly recovered. Thus, for the interval $[t_j + \alpha, t_j + 2\alpha)$, $\boldsymbol{e}_t$ again the expression given above but with $\hat{\boldsymbol{P}}_{j,0}$ replaced by $\hat{\boldsymbol{P}}_{j,1}$. Thus, for the second PCA update at $t = t_j + 2\alpha - 1$, $\|\boldsymbol{M}_{s,t}\boldsymbol{P}\|_2 \leq 1.2\mathrm{SE}(\hat{\boldsymbol{P}}_{j,1}, \boldsymbol{P}_j) \leq 1.2q_0/4 := q_1$. Apply the PCA-SDDN result with $q = q_1$ and with $\epsilon_{\mathrm{SE}} = q_1/4 = 1.2q_0/4^2$ to show that $\mathrm{SE}(\hat{\boldsymbol{P}}_{j,2}, \boldsymbol{P}_j) \leq q_1/4 = 1.2q_0/4^2$. Repeating this process, after $K$ updates, $\mathrm{SE}(\hat{\boldsymbol{P}}_{j,K}, \boldsymbol{P}_j) \leq (1.2/4)^{K-1}0.25q_0$. By picking $K$ large enough, we ensure that this bound is below $\epsilon$. We set the final estimate $\hat{\boldsymbol{P}}_j := \hat{\boldsymbol{P}}_{j,K}$. This serves as the starting point for estimating the next change.

In the above discussion, to explain things simply, in each subspace update, we used simple SVD. However, if we assume that only one (or only a few) direction(s) change at each subspace update time as was done in [23], one can get an improved result by first accurately estimating the new direction(s) that got added to the subspace by solving a problem of PCA with partial subspace knowledge repeated $K$ times. Finally, when $q$ is small enough (is below $2\epsilon$), one can re-estimate the entire subspace by solving a standard PCA problem. This is done to delete the removed direction(s) from the subspace estimate. For the former problem, we replace standard SVD by a projection-SVD step. To analyze it, we develop a modification of the ideas from this work to prove a result for PCA with partial subspace knowledge when noise is data dependent, see [23, Theorem 6.4]. For the deletion, step we use the PCA-SDDN result given above.
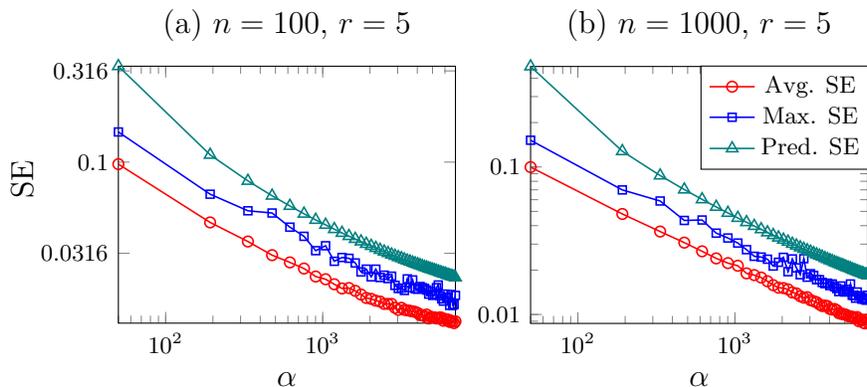
## 5   Numerical Experiments



Figure 1: Numerically computed mean and maximum values of $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P})$ (over 100 trials) and its bound from Theorem 2.13. The bound used $c = 1$.

In our first experiment we numerically demonstrate the tightness of the bound of Theorem 2.13 by plotting the numerically computed subspace error and the bound suggested by the theorem. In the expressions for $d(\alpha)$ and $d_{\mathrm{denom}}(\alpha)$ in the bound, there is an unspecified constant $c$. We set $c = 1$ while

plotting the bound. We generated the data as $\boldsymbol{\ell}_t = \boldsymbol{P}\boldsymbol{a}_t$, where $\boldsymbol{P}$ was generated by ortho-normalizing the columns of an $n \times r$ matrix with independent identically distributed (iid) standard Gaussian entries. We generated the coefficients $(\boldsymbol{a}_t)_i$ as iid $uniform(-6, 6)$. With this, $\lambda^+ = \lambda^- = 12$ and $f = 1$. We generated the uncorrelated noise as $\boldsymbol{v}_t = \boldsymbol{B}\boldsymbol{c}_t$ where $\boldsymbol{B}$ is generated by orthonormalizing the columns of an $n \times r_v$ matrix with iid standard Gaussian entries and $(\boldsymbol{c}_t) \sim unif(-q_i, q_i)$ with $q_i = 1.1 - 0.1i/r_v$. The data-dependent noise was generated as $\boldsymbol{w}_t = \boldsymbol{I}_{\mathcal{T}_t} \boldsymbol{M}_{s,t} \frac{q}{\|\boldsymbol{M}_{s,t}\boldsymbol{P}\|}\boldsymbol{\ell}_t$ and each entry of $\boldsymbol{M}_{s,t}$ was generated independently as the absolute value of a standard Gaussian r.v. (taking the absolute value ensures that $\mathbb{E}[\boldsymbol{M}_{s,t}] \neq 0$). Further, $\mathcal{T}_t$ was generated to follow [23, Model D.24] with $s = 5$, $\rho = 1$ and $b_0 = 0.05$ (simulates a 1D moving object that moves every so often). We set $\boldsymbol{y}_t = \boldsymbol{\ell}_t + \boldsymbol{w}_t + \boldsymbol{v}_t$. We used $r_v = r$, $q = 0.001$. From the support change model, $b = b_0 = 0.05$. We varied $\alpha$ in the range of $[29, 7000]$ and computed $\hat{\boldsymbol{P}}$ and $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P})$ for each value of $\alpha$. Fig 1b used $n = 1000$ and $r = 10$ while Fig. 1a used $n = 100$, $r = 5$. We show the mean and maximum values of the numerically computed $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P})$, and the bound predicted by Theorem 2.13, as a function of $\alpha$ in Fig. 1. The mean and max are computed over 100 Monte Carlo trials. Notice that the bound appears quite tight in both figures.
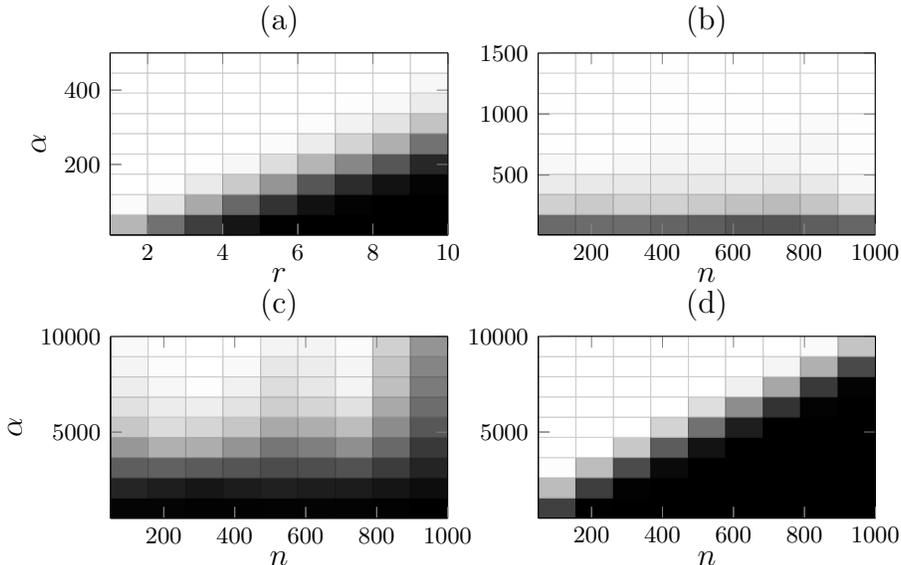


Figure 2: The grey scale intensity represents the numerically computed probability (fraction of times) that $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \le \epsilon$ for the $\epsilon$ value given in the text. Black is zero and white is one. Fig. 2a: displays the probability for $r$ versus $\alpha$ for $n = 100$, $r_v = r$ and bounded data and noise model; Fig. 2b: $n$ versus $\alpha$ for $r_v = r = 1$ and bounded data and noise; Fig. 2c: $n$ versus $\alpha$ for $r_v = r = 1$ and Gaussian data and noise; Fig. 2d: $n$ versus $\alpha$ for $r = 1$, $r_v = n$ and Gaussian.

In our second experiment, we use Monte Carlo to estimate the probability of $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \le \epsilon$ for various values of $r$ and $\alpha$ with $n$ fixed or various values

of $n$ and $\alpha$ with $r$ fixed. All our estimates used 100 Monte Carlo trails. The probability is displayed in Fig. 2 as a grey scale intensity with black denoting zero and white denoting one. This helps to numerically compute the value of $\alpha$ needed for a given $n, r$ to ensure that the probability is close to one (smallest $\alpha$ for which the color is white). Fig. 2a varies $r$ and $\alpha$ for $n = 100$. All other parameters were the same as in the first experiment. So $n = 100$, $b = 0.05$, $q = 0.001$, $f = 1$, $\lambda^- = 12$, $r_v = r$, $\lambda_v^+ = 1.1$. To generate the plot we used

$$\epsilon = 1.5 \left( \sqrt{b}(2q + q^2)f + \frac{\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}/\lambda^-}{1 - \frac{\lambda_{v,\mathrm{rest}} + -\lambda_{v,\boldsymbol{P}}^-}{\lambda^-}} \right).$$ As can be seen, the dependence of

$\alpha$ on $r$ is linear. This matches what our guarantees claim about the sample complexity: $\alpha$ needs to be $C \max(r_v, r) \log n$. Here $r_v = r$.

In the other three sub-figures we fix $r$ and $r_v$ and evaluate the dependence of $\alpha$ on $n$ in various settings. In Fig. 2b, we set $r_v = r$, $r = 1$, and other parameters were as above. Thus both true data and noise are bounded. We display the numerically estimated probability of $\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \le \epsilon$ for various values of $n$ and $\alpha$. Recall that for bounded data and noise, the required $\alpha$ is proportional to $\max(r_v, r) \log n$, i.e., it depends logarithmically on $n$. Logarithmic variation is hard to observe numerically unless a very large range of $n$ is used. This is also what is seen from Fig. 2b. For the range of values of $n$, the required $\alpha$ seems nearly constant.

For Fig. 2c, we replaced the boundedness assumption by a Gaussian assumption on data and noise. We still generated $\boldsymbol{v}_t = \boldsymbol{B}\boldsymbol{c}_t$ where $\boldsymbol{B}$ is an $n \times r_v$ matrix generated as before and we set $r_v = r$. But now we generated $(\boldsymbol{c}_t)_i \overset{\mathrm{iid}}{\sim} \mathcal{N}(0, q_i^2)$ with $q_i = 0.9 - 0.4i/r_v$. We generated $\boldsymbol{\ell}_t = \boldsymbol{P}\boldsymbol{a}_t$ where $\boldsymbol{a}_t \overset{\mathrm{iid}}{\sim} \mathcal{N}(0, 100)$. Here $\mathcal{N}(0, \sigma^2)$ refers to a zero Gaussian distribution with variance $\sigma^2$. Thus $\boldsymbol{\Lambda} = 100$ and $f = 1$. Everything else was the same as in the first experiment, thus $\boldsymbol{w}_t = \boldsymbol{I}_{\mathcal{T}_t} \boldsymbol{M}_{s,t} \boldsymbol{\ell}_t$ with $\mathcal{T}_t$ and $\boldsymbol{M}_{s,t}$ were generated as described earlier. Thus, $b = 0.05$, $q = 0.001$. Also, $\lambda^+ = 100 = \lambda^-$ and $f = 1$ and $\lambda_v^+ = 0.9$. We fixed $r_v = r = 1$ and vary $n$ and $\alpha$ as was also done in Fig. 2b. Notice now that the required $\alpha$ to achieve high-enough (white) probability of success does increase with $n$. It is hard to say though whether the dependence is indeed linear as predicted by our theorem. As we pointed out earlier in Remark 2.11, for this setting since $r_v = r \ll n$, it may be possible to tighten the required sample complexity lower bound.

Finally, for Fig. 2d, we generated data exactly as for Fig. 2c but with $r_v = n$ (instead of $r_v = r$). As can be seen, now the required sample complexity does indeed increase linearly with $n$. This matches what is predicted by our main result for the $r_v = n$ case. Notice that in the last two sub-figures where we used Gaussian noise and data, we have increased signal power as compared to the bounded case. If this was not increased, the required $\alpha$ to achieve large enough probability of success would be very large and would lead to very slow computations. All experiments used the MATLAB command svds for computing $\hat{\boldsymbol{P}}$. All codes are available at `https://github.com/praneethmurthy/correlated-pca`.

# 6    Conclusions and Future Work

In this work, we studied the PCA problem when the noise can be non-isotropic and/or data-dependent, and as a result, in general the data and noise are correlated. We obtained guarantees under both a bounded-ness assumption and a sub-Gaussian assumption on the data and noise. When the uncorrelated noise has effective dimension $O(r)$, under the bounded-ness assumption, a simple assumption on data-noise correlation, and a bound on the ratio between noise power and minimum signal space eigenvalue, we showed that the required sample complexity for PCA is near optimal. Under the sub-Gaussian assumption, the required sample complexity as predicted by our results increases to $O(n)$ which is comparable to what existing results for isotropic Gaussian noise also need. However, as noted in Remark 2.11, in the setting where the sub-Gaussian noise has effective dimension $r_v \ll n$, it should be possible to tighten this.

The result given here assumes that the $\boldsymbol{\ell}_t$'s are mutually independent random variables. Mutual independence can be replaced by an autoregressive (AR) model on the $\boldsymbol{\ell}_t$'s. As long as the AR parameter is not too large, it should be possible to get a result very similar to the one given in this work using the matrix Freedman's inequality [25] or a little weaker than the one given here using matrix Azuma [15]. The latter would generalize the approach developed in [9] to for analyzing the subspace update step of ReProCS under an AR model on the $\boldsymbol{\ell}_t$'s.

In ongoing work, we are studying the problem of PCA in data-dependent noise when partial knowledge of the subspace is available and its implications for the subspace update step of ReProCS [23]. A useful open question for future work is how to analyze algorithms for streaming PCA, e.g., the block-stochastic power method, in the data-dependent noise setting. This was studied in [4] under the spiked covariance model, or in [5] for an arbitrary observed data covariance matrix, but for $r = 1$ dimensional PCA.

# A    Proof of Main Result: Bounded Case

*Proof of Theorem 2.13 under Assumption 2.4.* Apply the Davis-Kahan $\sin\theta$ theorem [14] summarized in Lemma 2.8 with

$$\boldsymbol{D}_0 = \frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t \boldsymbol{\ell}_t' + \boldsymbol{P}\boldsymbol{P}'\boldsymbol{\Sigma}_v \boldsymbol{P}\boldsymbol{P}' = \boldsymbol{P}(\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t' + \boldsymbol{P}'\boldsymbol{\Sigma}_v \boldsymbol{P})\boldsymbol{P}'.$$

Observe that $\lambda_{r+1}(\boldsymbol{D}_0) = 0$ and, using Weyl and $\mathbb{E}[\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t'] = \boldsymbol{\Lambda}$,

$$\lambda_r(\boldsymbol{D}_0) \geq \lambda_{\min}(\boldsymbol{P}'\boldsymbol{\Sigma}_v \boldsymbol{P}) + \lambda_{\min}(\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t') \geq \lambda_{v,\boldsymbol{P}}^- + \lambda^- - \|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda}\|_2.$$

Thus, rewriting $\boldsymbol{D} - \boldsymbol{D}_0$ as $\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0] + (\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0])$,

$$\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \leq \frac{\|(\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0])\boldsymbol{P}\|_2 + \|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2}{\lambda^- + \lambda_{v,\boldsymbol{P}}^- - \|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda}\|_2 - \lambda_{\max}(\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]) - \|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2}$$

Observe that

$$\boldsymbol{D} - \boldsymbol{D}_0 = \frac{1}{\alpha}\sum_t \boldsymbol{v}_t \boldsymbol{v}_t' + \frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t \boldsymbol{v}_t' + \frac{1}{\alpha}\sum_t \boldsymbol{v}_t \boldsymbol{\ell}_t'$$
$$+ \frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t \boldsymbol{w}_t' + \frac{1}{\alpha}\sum_t \boldsymbol{w}_t \boldsymbol{\ell}_t' + \frac{1}{\alpha}\sum_t \boldsymbol{w}_t \boldsymbol{w}_t' - \boldsymbol{P}\boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}\boldsymbol{P}'$$

Since $\boldsymbol{v}_t$ is uncorrelated with $\boldsymbol{\ell}_t$, the expected value of the second and third terms is zero. Thus,

$$\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0] = (\boldsymbol{\Sigma}_v - \boldsymbol{P}\boldsymbol{P}'\boldsymbol{\Sigma}_v\boldsymbol{P}\boldsymbol{P}') + \frac{1}{\alpha}\sum_t \boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_t' + \frac{1}{\alpha}\sum_t \boldsymbol{M}_t\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}' + \frac{1}{\alpha}\sum_t \boldsymbol{M}_t\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_t'$$
$$(10)$$

The third term on the RHS is a transpose of the second one. For the second term, write $\boldsymbol{M}_t = \boldsymbol{M}_{2,t}\boldsymbol{M}_{1,t}$ and then use Cauchy Schwartz for sums of matrices (see e.g. [9, Lemma A.6]) to conclude that

$$\|\frac{1}{\alpha}\sum_t \boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_{1,t}'\boldsymbol{M}_{2,t}'\|_2^2 \le \|\frac{1}{\alpha}\sum_t \boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_{1,t}'\boldsymbol{M}_{1,t}\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\|_2 \|\frac{1}{\alpha}\sum_t \boldsymbol{M}_{2,t}\boldsymbol{M}_{2,t}'\|_2$$
$$\le \max_t \|\boldsymbol{M}_{1,t}\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\|_2^2 \, b \le (q\lambda^+)^2 b. \qquad (11)$$

The second inequality used (5) and the third inequality used $\|\boldsymbol{M}_{1,t}\boldsymbol{P}\| \le q$. The fourth term is handled similarly:

$$\|\frac{1}{\alpha}\sum_t \boldsymbol{M}_t\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_{1,t}'\boldsymbol{M}_{2,t}'\|_2^2 \le \|\frac{1}{\alpha}\sum_t \boldsymbol{M}_t\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_{1,t}'\boldsymbol{M}_{1,t}\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_t'\|_2 \|\frac{1}{\alpha}\sum_t \boldsymbol{M}_{2,t}\boldsymbol{M}_{2,t}'\|_2$$
$$\le \max_t \|\boldsymbol{M}_t\boldsymbol{P}\boldsymbol{\Lambda}\boldsymbol{P}'\boldsymbol{M}_{1,t}'\|_2^2 \, b \le (q^2\lambda^+)^2 b. \quad (12)$$

Since $(\boldsymbol{P}\boldsymbol{P}' + \boldsymbol{P}_\perp \boldsymbol{P}_\perp') = \boldsymbol{I}$, we can always rewrite $\boldsymbol{\Sigma}_v = (\boldsymbol{P}\boldsymbol{P}' + \boldsymbol{P}_\perp \boldsymbol{P}_\perp')\boldsymbol{\Sigma}_v(\boldsymbol{P}\boldsymbol{P}' + \boldsymbol{P}_\perp \boldsymbol{P}_\perp')$. Using this and the last three equations above,

$$\|(\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0])\boldsymbol{P}\|_2 \le \|\boldsymbol{P}_\perp'\boldsymbol{\Sigma}_v\boldsymbol{P}\|_2 + 2\sqrt{b}q\lambda^+ + \sqrt{b}q^2\lambda^+ = \lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp} + \sqrt{b}(2q + q^2)\lambda^+$$

and

$$\lambda_{\max}(\mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]) \le \lambda_{v,\mathrm{rest}}^+ + \sqrt{b}(2q + q^2)\lambda^+.$$

Thus, by Lemma 2.8,

$$\mathrm{SE}(\hat{\boldsymbol{P}}, \boldsymbol{P}) \le \frac{\frac{\lambda_{v,\boldsymbol{P},\boldsymbol{P}_\perp}}{\lambda^-} + \sqrt{b}(2q + q^2)f + \frac{\|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2}{\lambda^-}}{1 - (\frac{\lambda_{v,\mathrm{rest}}^+ - \lambda_{v,\boldsymbol{P}}^-}{\lambda^-} + \sqrt{b}(2q + q^2)f) - \frac{\|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda}\|_2}{\lambda^-} - \frac{\|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2}{\lambda^-}}$$
$$(13)$$

To bound $\|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2$ and $\|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda}\|_2$, we use concentration bounds from the following lemma which we prove in Appendix B.

23

**Lemma A.20.** *With probability at least $1 - 10n^{-10}$, if $\alpha^3 > \max(r_v, r) \log n$, then,*

$$\left\| \frac{1}{\alpha} \sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda} \right\| \leq c\eta f \sqrt{\frac{r + \log n}{\alpha}} \lambda^-,$$

$$\left\| \frac{1}{\alpha} \sum_t \boldsymbol{\ell}_t \boldsymbol{w}_t' - \frac{1}{\alpha} \mathbb{E}[\sum_t \boldsymbol{\ell}_t \boldsymbol{w}_t'] \right\|_2 \leq c\sqrt{\eta} q f \sqrt{\frac{r \log n}{\alpha}} \lambda^-,$$

$$\left\| \frac{1}{\alpha} \sum_t \boldsymbol{w}_t \boldsymbol{w}_t' - \frac{1}{\alpha} \mathbb{E}[\sum_t \boldsymbol{w}_t \boldsymbol{w}_t'] \right\|_2 \leq c\sqrt{\eta} q^2 f \sqrt{\frac{r \log n}{\alpha}} \lambda^-,$$

$$\left\| \frac{1}{\alpha} \sum_t \boldsymbol{\ell}_t \boldsymbol{v}_t' \right\|_2 \leq c\sqrt{\eta} \sqrt{\frac{\lambda_v^+}{\lambda^-}} f \sqrt{\frac{\max(r_v, r) \log n}{\alpha}} \lambda^-,$$

$$\left\| \frac{1}{\alpha} \sum_t \boldsymbol{v}_t \boldsymbol{v}_t' - \frac{1}{\alpha} \mathbb{E}[\sum_t \boldsymbol{v}_t \boldsymbol{v}_t'] \right\|_2 \leq c\sqrt{\eta} \frac{\lambda_v^+}{\lambda^-} \sqrt{\frac{r_v \log n}{\alpha}} \lambda^-.$$

Thus, w.p. $\geq 1 - 10n^{-10}$, $\left\| \frac{1}{\alpha} \sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda} \right\| \leq c\eta f \sqrt{\frac{r + \log n}{\alpha}} \lambda^-$ and $\|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2 \leq d(\alpha)_0 \lambda^-$ where

$$d(\alpha)_0 := c\sqrt{\eta} \max\left( q f \sqrt{\frac{r \log n}{\alpha}}, \sqrt{\frac{\lambda_v^+}{\lambda^-}} f \sqrt{\frac{\max(r_v, r) \log n}{\alpha}}, \frac{\lambda_v^+}{\lambda^-} \sqrt{\frac{r_v \log n}{\alpha}} \right)$$
(14)

Clearly, $d(\alpha)_0 \leq d(\alpha)$ defined in the statement of Theorem 2.13. Combining these bounds with (13) we get our final result. ⊠

# B  Proof of Concentration Bounds

*Proof of Lemma A.20.*

$\boldsymbol{a}_t \boldsymbol{a}_t'$ *term.*  Using Vershynin's sub-Gaussian result (Theorem 5.39 of [13]) applied to $\frac{1}{\alpha} \sum_t \boldsymbol{a}_t \boldsymbol{a}_t'$, and using the fact that the $\boldsymbol{a}_t$'s are $r$-length independent sub-Gaussian vectors with sub-Gaussian norm bounded by $\sqrt{\eta \lambda^+}$, we get the following: with probability at least $1 - 2\exp\left(r \log 9 - \alpha \frac{c(\epsilon_1 \lambda^-)^2}{(4\eta\lambda^+)^2}\right) = 1 - 2\exp\left(r \log 9 - \alpha \frac{c\epsilon_1^2}{16\eta^2 f^2}\right)$,

$$\|\frac{1}{\alpha} \sum_t \boldsymbol{a}_t \boldsymbol{a}_t' - \boldsymbol{\Lambda}\|_2 \leq \epsilon_1 \lambda^-$$

Set $\epsilon_1 = c\eta f \sqrt{\frac{r + 11 \log n}{\alpha}}$. Then, the above event holds w.p. at least $1 - 2n^{-10}$

$\boldsymbol{\ell}_t \boldsymbol{w}_t'$ *term.*  This and all other items use Matrix Bernstein for rectangular matrices, Theorem 1.6 of [15]. This says the following. For a finite sequence of

$d_1 \times d_2$ zero mean independent matrices $\boldsymbol{Z}_k$ with

$$\|\boldsymbol{Z}_k\|_2 \leq R, \text{ and} \max(\| \sum_k \mathbb{E}[\boldsymbol{Z}_k'\boldsymbol{Z}_k]\|_2, \| \sum_k \mathbb{E}[\boldsymbol{Z}_k\boldsymbol{Z}_k']\|_2) \leq \sigma^2,$$

we have $\Pr(\| \sum_k \boldsymbol{Z}_k\|_2 \geq s) \leq (d_1 + d_2) \exp\left(-\frac{s^2/2}{\sigma^2 + Rs/3}\right)$.

Let $\boldsymbol{Z}_t := \boldsymbol{\ell}_t \boldsymbol{w}_t'$. We apply this result to $\tilde{\boldsymbol{Z}}_t := \boldsymbol{Z}_t - \mathbb{E}[\boldsymbol{Z}_t]$ with $s = \epsilon\alpha$. To get the values of $R$ and $\sigma^2$ in a simple fashion, we use the facts that (i) if $\|\boldsymbol{Z}_t\|_2 \leq R_1$, then $\|\tilde{\boldsymbol{Z}}_t\| \leq 2R_1$; and (ii) $\sum_t \mathbb{E}[\tilde{\boldsymbol{Z}}_t\tilde{\boldsymbol{Z}}_t'] \preccurlyeq \sum_t \mathbb{E}[\boldsymbol{Z}_t\boldsymbol{Z}_t']$. Thus, we can set $R$ to two times the bound on $\|\boldsymbol{Z}_t\|_2$ and we can set $\sigma^2$ as the maximum of the bounds on $\| \sum_t \mathbb{E}[\boldsymbol{Z}_t\boldsymbol{Z}_t']\|_2$ and $\| \sum_t \mathbb{E}[\boldsymbol{Z}_t'\boldsymbol{Z}_t]\|_2$.

It is easy to see that $R = 2\sqrt{\eta r\lambda^+}\sqrt{\eta rq^2\lambda^+} = 2\eta rq\lambda^+$. To get $\sigma^2$, observe that

$$\left\| \sum_t \mathbb{E}[\boldsymbol{w}_t\boldsymbol{\ell}_t'\boldsymbol{\ell}_t\boldsymbol{w}_t'] \right\|_2 \leq \alpha(\max_{\boldsymbol{\ell}_t}\|\boldsymbol{\ell}_t\|^2) \cdot \|\mathbb{E}[\boldsymbol{w}_t\boldsymbol{w}_t']\|$$
$$\leq \alpha\eta r\lambda^+ \cdot q^2\lambda^+ = \alpha\eta rq^2(\lambda^+)^2.$$

Repeating the above steps, we get the same bound on $\| \sum_t \mathbb{E}[\boldsymbol{Z}_t\boldsymbol{Z}_t']\|_2$. Thus, $\sigma^2 = \alpha\eta rq^2(\lambda^+)^2$.

Thus, we conclude that,

$$\left\| \sum_t \boldsymbol{\ell}_t\boldsymbol{w}_t' - \mathbb{E}[\sum_t \boldsymbol{\ell}_t\boldsymbol{w}_t'] \right\|_2 \geq \epsilon\alpha \tag{15}$$

w.p. at most $2n\exp\left(-\frac{\epsilon^2\alpha/2}{\eta rq^2(\lambda^+)^2 + (\eta rq\lambda^+ \epsilon/3\alpha)}\right)$.

Set $\epsilon = \epsilon_0\lambda^-$ with $\epsilon_0 = c\sqrt{\eta}qf\sqrt{\frac{r\log n}{\alpha}}$. If $\alpha^3 > \eta(r\log n)$, then clearly, $R\epsilon \leq \sigma^2$. With this $\epsilon$, if $\alpha^3 > (r\log n)$, (15) holds w.p. at least $1 - 2n^{-10}$.

$\boldsymbol{w}_t\boldsymbol{w}_t'$ *term.* We again apply matrix Bernstein and proceed as above. In this case, $R = 2\eta rq^2\lambda^+$ and $\sigma^2 = \alpha\sigma_1^2$, $\sigma_1^2 = \eta rq^4(\lambda^+)^2$. Thus $R < \sigma^2/2q^2$. Set $\epsilon = \epsilon_2\lambda^-$ with $\epsilon_2 = c\sqrt{\eta}q^2f\sqrt{\frac{r\log n}{\alpha}} < 1$. Then, can again show that if $\alpha^3 > (r\log n)$, the probability of the bad event is bounded by $2n^{-10}$.

$\boldsymbol{\ell}_t\boldsymbol{v}_t'$ *and* $\boldsymbol{v}_t\boldsymbol{v}_t'$ *terms.* Apply matrix Bernstein as done above. $\boxtimes$

# C   Proof of Main Result: Sub-Gaussian Case

*Proof of Theorem 2.13 under Assumption 2.5.* The only thing that changes in this case is the $\epsilon$ values used for the various terms in Lemma A.20. These change because we cannot apply matrix Bernstein now. Instead, for all terms, we apply the following simple modification of Vershynin's sub-Gaussian result [13, Theorem 5.39]. This lemma follows using exactly the proof approach of [13, Theorem 5.39] but with the following change. If two r.v.s $x$, $y$ are sub-Gaussian with sub-Gaussian norms $k_x$, $k_y$ respectively, then the r.v. $z := xy$ is sub-exponential with

sub-exponential norm bounded by $ck_x k_y$. This fact itself follows by Cauchy-Schwartz and the definitions of the sub-Gaussian and sub-exponential norms:
$\mathbb{E}[|xy|^p]^{1/p} \leq ((\mathbb{E}[|x|^{2p}]\mathbb{E}[|y|^{2p}])^{1/2})^{1/p} = (\mathbb{E}[|x|^{2p}]\mathbb{E}[|y|^{2p}])^{1/2p} \leq ((\sqrt{2p}k_x)^{2p}(\sqrt{2p}k_y)^{2p})^{1/2p} = (\sqrt{2p}k_x)(\sqrt{2p}k_y) = p(2k_x k_y)$. The first inequality used Cauchy-Schwartz, the second inequality used the definition of sub-Gaussian norm [13, Sec. 5.2].

**Lemma C.21.** *Let $\boldsymbol{x}_i$, $i = 1, 2, \ldots, N$, and $\boldsymbol{y}_i, i = 1, 2, \ldots, N$ be zero mean sub-Gaussian random vectors with sub-Gaussian norms bounded by $K_x$ and $K_y$ respectively. Each $\boldsymbol{x}_i, \boldsymbol{y}_i$ is in $\mathbb{R}^n$. Also $\{\boldsymbol{x}_i, \boldsymbol{y}_i\}$, $i = 1, 2, \ldots, N$ are mutually independent. Then,*

$$\Pr\left(\left\|\frac{1}{N}\sum_i \boldsymbol{x}_i {\boldsymbol{y}_i}' - \mathbb{E}\left[\frac{1}{N}\sum_i \boldsymbol{x}_i {\boldsymbol{y}_i}'\right]\right\| \leq t\right) \geq 1 - 2\exp\left(n\log 9 - \frac{t^2 N}{16c(K_x K_y)}\right)$$

Using the above lemma, we can conclude the following.

1. Bounding $\left\|\frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t {\boldsymbol{w}_t}' - \frac{1}{\alpha}\mathbb{E}[\sum_t \boldsymbol{\ell}_t {\boldsymbol{w}_t}']\right\|_2$: Apply Lemma C.21 with $\boldsymbol{x}_t = \boldsymbol{\ell}_t$, $\boldsymbol{y}_t = \boldsymbol{w}_t$, $N \equiv \alpha$. Then $K_x = c\sqrt{\lambda^+}$ and $K_y = c\sqrt{q^2\lambda^+}$, thus $K_x K_y = cq\lambda^+$. Set $t = \epsilon_0 \lambda^-$ with $\epsilon_0 = cf\sqrt{\frac{n}{\alpha}}$. Then

$$\Pr\left(\left\|\frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t {\boldsymbol{w}_t}' - \frac{1}{\alpha}\mathbb{E}[\sum_t \boldsymbol{\ell}_t {\boldsymbol{w}_t}']\right\|_2 \leq \epsilon_0 \lambda^-\right) \geq 1 - 2\exp(-cn)$$

2. Bounding $\left\|\frac{1}{\alpha}\sum_t \boldsymbol{w}_t {\boldsymbol{w}_t}' - \frac{1}{\alpha}\mathbb{E}[\sum_t \boldsymbol{w}_t {\boldsymbol{w}_t}']\right\|_2$: Proceed as above. Use $\epsilon = \epsilon_0 \lambda^-$.

3. Bounding $\left\|\frac{1}{\alpha}\sum_t \boldsymbol{a}_t {\boldsymbol{a}_t}' - \boldsymbol{\Lambda}\right\|_2$: Proceed as above. Use $\epsilon = \epsilon_0 \lambda^-$.

4. Bounding $\left\|\frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t {\boldsymbol{v}_t}' - \frac{1}{\alpha}\mathbb{E}[\sum_t \boldsymbol{\ell}_t {\boldsymbol{v}_t}']\right\|_2$: Apply Lemma C.21 with $\boldsymbol{x}_t = \boldsymbol{\ell}_t$, $\boldsymbol{y}_t = \boldsymbol{v}_t$, $N \equiv \alpha$. Then $K_x = c\sqrt{\lambda^+}$ and $K_y = c\sqrt{\lambda_v^+}$, thus $K_x K_y = c\sqrt{\lambda^+ \lambda_v^+}$ Set $t = \epsilon_{0,v}\lambda^-$ with $\epsilon_{0,v} = c\sqrt{\frac{\lambda_v^+}{\lambda^-}}f\sqrt{\frac{n}{\alpha}}$. Then,

$$\Pr\left(\left\|\frac{1}{\alpha}\sum_t \boldsymbol{\ell}_t {\boldsymbol{v}_t}' - \frac{1}{\alpha}\mathbb{E}[\sum_t \boldsymbol{\ell}_t {\boldsymbol{v}_t}']\right\|_2 \leq \epsilon_{0,v}\lambda^-\right) \geq 1 - 2\exp(-cn)$$

5. Bounding $\left\|\frac{1}{\alpha}\sum_t \boldsymbol{v}_t {\boldsymbol{v}_t}' - \frac{1}{\alpha}\mathbb{E}[\sum_t \boldsymbol{v}_t {\boldsymbol{v}_t}']\right\|_2$: Apply [13, Theorem 5.39] and proceed as in the previous part. Use $\epsilon = \epsilon_{1,v}\lambda^-$ with $\epsilon_{1,v} = \frac{\lambda_v^+}{\lambda^-}\sqrt{\frac{n}{\alpha}}$.

Combining the above bounds (except the third one), we conclude that, w.p. $\geq 1 - 10\exp(-cn)$,

$$\|\boldsymbol{D} - \boldsymbol{D}_0 - \mathbb{E}[\boldsymbol{D} - \boldsymbol{D}_0]\|_2 \leq 5d(\alpha)_{sG} := c\max\left(\frac{\lambda_v^+}{\lambda^-}, f\right)\sqrt{\frac{n}{\alpha}} \qquad (16)$$

Everything else remains the same. $\boxtimes$

# References

[1] D. Hong, L. Balzano, and J. A. Fessler, "Towards a theoretical analysis of pca for heteroscedastic data," in *Allerton Conf. Comm. Control Comput.*, 2016.

[2] B. Nadler, "Finite sample approximation results for principal component analysis: A matrix perturbation approach," *Ann. Statist.*, vol. 36, no. 6, 2008.

[3] Iain M Johnstone, "On the distribution of the largest eigenvalue in principal components analysis," *Ann. Statist.*, pp. 295–327, 2001.

[4] I. Mitliagkas, C. Caramanis, and P. Jain, "Memory limited, streaming pca," in *Adv. Neural Info. Proc. Sys. (NIPS)*, 2013, pp. 2886–2894.

[5] Prateek Jain, Chi Jin, Sham M Kakade, Praneeth Netrapalli, and Aaron Sidford, "Streaming pca: Matching matrix bernstein and near-optimal finite sample guarantees for ojas algorithm," in *29th Annual Conference on Learning Theory*, 2016, pp. 1147–1164.

[6] N. Vaswani and H. Guo, "Correlated-pca: Principal components' analysis when data and noise are correlated," in *Adv. Neural Info. Proc. Sys. (NIPS)*, 2016.

[7] Jussi Gillberg, Pekka Marttinen, Matti Pirinen, Antti J Kangas, Pasi Soininen, Mehreen Ali, Aki S Havulinna, Marjo-Riitta Järvelin, Mika Ala-Korpela, and Samuel Kaski, "Multiple output regression with latent noise," *J. Mach. Learn. Res.*, 2016.

[8] C. Qiu, N. Vaswani, B. Lois, and L. Hogben, "Recursive robust pca or recursive sparse recovery in large but structured noise," *IEEE Trans. Info. Th.*, pp. 5007–5039, August 2014.

[9] J. Zhan, B. Lois, H. Guo, and N. Vaswani, "Online (and Offline) Robust PCA: Novel Algorithms and Performance Guarantees," in *Intnl. Conf. Artif. Intell. Stat. (AISTATS)*, 2016, long version: ArXiv: 1601.07985.

[10] A. Beck and M. Teboulle, "Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems," *IEEE Trans. Image Process.*, vol. 18, no. 11, pp. 2419–2434, 2009.

[11] Ren-Cang Li, "Relative perturbation theory: Ii. eigenspace and singular subspace variations," *SIAM J. Matrix Anal. Appl.*, vol. 20, no. 2, pp. 471–492, 1998.

[12] R. Vershynin, "How close is the sample covariance matrix to the actual covariance matrix?," *J. Theoret. Probab.*, pp. 1–32, 2010.

[13] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," *Compressed sensing*, pp. 210–268, 2012.

[14] C. Davis and W. M. Kahan, "The rotation of eigenvectors by a perturbation. iii," *SIAM J. Numer. Anal.*, vol. 7, pp. 1–46, Mar. 1970.

[15] J. A. Tropp, "User-friendly tail bounds for sums of random matrices," *Found. Comput. Math.*, vol. 12, no. 4, 2012.

[16] P. Netrapalli, P. Jain, and S. Sanghavi, "Phase retrieval using alternating minimization," in *Adv. Neural Info. Proc. Sys. (NIPS)*, 2013, pp. 2796–2804.

[17] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?," *J. ACM*, vol. 58, no. 3, 2011.

[18] P. Netrapalli, U N Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain, "Non-convex robust pca," in *Neural Info. Proc. Sys. (NIPS)*, 2014.

[19] Xinyang Yi, Dohyung Park, Yudong Chen, and Constantine Caramanis, "Fast algorithms for robust pca via gradient descent," in *Neural Info. Proc. Sys. (NIPS)*, 2016.

[20] Yeshwanth Cherapanamjeri, Kartik Gupta, and Prateek Jain, "Nearly-optimal robust matrix completion," *arXiv preprint arXiv:1606.07315*, 2016.

[21] H. Xu, C. Caramanis, and S. Sanghavi, "Robust pca via outlier pursuit," *IEEE Trans. Inform. Theory*, vol. 58, no. 5, May 2012.

[22] P. Netrapalli, P. Jain, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Symposium on Theory of Computing (STOC)*, 2013.

[23] P. Narayanamurthy and N. Vaswani, "New Results for Provable Dynamic Robust PCA," *arXiv:1705.08948*, 2017.

[24] E. Candes, "The restricted isometry property and its implications for compressed sensing," *C. R. Math. Acad. Sci. Paris. Serie I*, pp. 589–592, 2008.

[25] J. A. Tropp, "Freedmans inequality for matrix martingales," *Electron. Commun. Probab*, vol. 16, pp. 262–270, 2011.