
Sparse Principal Component Analysis for High Dimensional Multivariate Time Series

Zhaoran Wang
Department of Electrical Engineering,
Princeton University

Fang Han
Department of Biostatistics,
Johns Hopkins University

Han Liu
Department of Operations Research
and Financial Engineering,
Princeton University

Abstract

We study sparse principal component analysis (sparse PCA) for high dimensional multivariate vector autoregressive (VAR) time series. By treating the transition matrix as a nuisance parameter, we show that sparse PCA can be directly applied on analyzing multivariate time series as if the data are i.i.d. generated. Under a double asymptotic framework in which both the length of the sample period T and dimensionality d of the time series can increase (with possibly $d \gg T$), we provide explicit rates of convergence of the angle between the estimated and population leading eigenvectors of the time series covariance matrix. Our results suggest that the spectral norm of the transition matrix plays a pivotal role in determining the final rates of convergence. Implications of such a general result is further illustrated using concrete examples. The results of this paper have impacts on different applications, including financial time series, biomedical imaging, and social media, etc.

1 Introduction

This paper considers sparse principal component analysis for weakly stationary vector autoregressive (VAR) time series (In this paper, we only consider VAR(1) model, i.e., the model with lag 1. We extend our results to VAR(p) in the longer version of this paper): Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \in \mathbb{R}^d$ be T observations from a time series $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_T$. Here we assume that each $\mathbf{X}_t \in \mathbb{R}^d$ is a d -dimensional random vector and

follows a VAR model

$$\mathbf{X}_{t+1} = \mathbf{A}\mathbf{X}_t + \mathbf{Z}_t, \text{ for } t = 1, 2, \dots, T-1, \quad (1.1)$$

where $\mathbf{A} \in \mathbb{R}^{d \times d}$ is called transition matrix and $\mathbf{Z}_1, \mathbf{Z}_2, \dots \stackrel{\text{i.i.d.}}{\sim} N_d(\mathbf{0}, \Psi)$ are independent colored Gaussian noise with covariance matrix Ψ . Since the process is weakly stationary, we denote Σ to be the covariance matrix of the time series,

$$\Sigma := \text{Var}(\mathbf{X}_1) = \dots = \text{Var}(\mathbf{X}_T). \quad (1.2)$$

Let $\mathbf{u}_1, \dots, \mathbf{u}_m$ be the top m leading eigenvectors of Σ . We want to find $\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_m$ which can estimate $\mathbf{u}_1, \dots, \mathbf{u}_m$ accurately. In the special case where \mathbf{A} is a zero matrix, this problem reduces to classical principal component analysis for i.i.d. Gaussian data. In more general settings where $\mathbf{A} \neq \mathbf{0}$, it is well known that to secure the weakly stationarity of the time series in (1.1), we must have the spectral norm of \mathbf{A} (i.e., the largest singular value of \mathbf{A}) smaller than 1.

In this paper we consider high dimensional time series under a double asymptotic framework, i.e., we allow the time series dimension d to scale with the length of the sample period T with possibly $d \gg T$. Compared with the classical asymptotic framework for time series in which only T increases while d remains fixed, such a theoretical framework better reflects the challenge in many real-world applications. For example, in fMRI image processing, the machine collects T scans of the human brain, each of which contains d voxels. In a typical setting, the number of scans T is around hundreds, while the number of voxels d could be tens of thousands. In another application on modeling social media stream, e.g., twitter data, we simultaneously monitor the number of tweets for d persons across T time units (e.g. hours). In a typical setting, T could be hundreds or thousands while d could be millions. Other applications include low-frequency stock data in which T represents the number of records of the closing price and d represents the number of stocks in the market.

Such a double asymptotic framework, though more re-

Appearing in Proceedings of the 16th International Conference on Artificial Intelligence and Statistics (AISTATS) 2013, Scottsdale, AZ, USA. Volume 31 of JMLR: W&CP 31. Copyright 2013 by the authors.

alistic, poses significant theoretical challenges. Even in a simplified setting where $\mathbf{A} = \mathbf{0}$, Johnstone and Lu (2009) show that the classical PCA is inconsistent under some conditions. In other words, letting the estimator $\hat{\mathbf{u}}_1$ be the leading eigenvector of the sample covariance matrix, the angle between \mathbf{u}_1 and $\hat{\mathbf{u}}_1$, denoted as $\angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)$, doesn't converge to 0 as T goes to infinity. To avoid such a curse of dimensionality, the population leading eigenvector \mathbf{u}_1 is in general assumed to be sparse. More specifically, let s be the number of nonzero elements in \mathbf{u}_1 , we assume $s \ll T$. With this sparsity assumption, different versions of sparse PCA have been proposed to handle i.i.d. Gaussian data (i.e., $\mathbf{A} = \mathbf{0}$): For example, greedy algorithms (d'Aspremont et al., 2008), lasso-type methods including SCoTLASS (Jolliffe et al., 2003), SPCA (Zou, 2006) and sPCA-rSVD (Shen and Huang, 2008), a number of power methods (Journée et al., 2010; Yuan, 2010; Ma, 2011), the biconvex algorithm PMD (Witten et al., 2009) and the semidefinite relaxation DSPCA (d'Aspremont et al., 2004). Sparse PCA has been widely used in finance (d'Aspremont et al., 2005), text mining (Zhang and El Ghaoui, 2011) and voting data analysis (Zhang et al., 2012).

One drawback for these existing sparse PCA theories is that they all assume the T observations $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ are independently and identically distributed. Such an assumption is obviously violated in most real-world applications. For example, in the fMRI imaging application we described before, the scans from two adjacent time points are obviously correlated. In the applications for stock price or twitter data, the existence of non-dependence is easily justified. Though some work exists for low-dimensional PCA on dependent data (Skinner et al., 1986), no such result exists for high dimensional settings. There are some related results on dependent data analysis in high dimensions (Loh and Wainwright, 2011; Fan et al., 2012), they are mainly for other learning methods. For example, Loh and Wainwright (2011) study the high dimensional regression for Gaussian data with missing values and dependent data. Very recently, Fan et al. (2012) analyze the penalized least square estimators, taking a weakly dependence structure, called α -mixing, of the noisy term into consideration.

In this paper, we study sparse PCA for weakly stationary VAR time series. By treating the transition matrix \mathbf{A} as a nuisance parameter, we directly apply sparse PCA on the multivariate time series $\mathbf{x}_1, \dots, \mathbf{x}_T$ to estimate the leading eigenvector \mathbf{u}_1 as if the data are i.i.d. generated. Let $\angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)$ be the angle between $\hat{\mathbf{u}}_1$ and \mathbf{u}_1 , and $\lambda_k(\boldsymbol{\Sigma})$ be the k -th largest value of $\boldsymbol{\Sigma}$, we provide an explicit rate of convergence for $\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)$, i.e., for some absolute constant C ,

$$|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)| \leq \frac{C}{\lambda_1(\boldsymbol{\Sigma}) - \lambda_2(\boldsymbol{\Sigma})} \sqrt{\frac{s \log d}{T} \left(\frac{\|\boldsymbol{\Sigma}\|_2}{1 - \|\mathbf{A}\|_2} \right)},$$

with high probability.

Our result allows the quantities $\lambda_1(\boldsymbol{\Sigma})$, $\lambda_2(\boldsymbol{\Sigma})$ and $\|\mathbf{A}\|_2$ all scale with d . Since $\boldsymbol{\Sigma}$ is jointly determined by the structure of the transition matrix \mathbf{A} and noise covariance matrix $\boldsymbol{\Psi}$, this result suggests that the interaction between the structure of \mathbf{A} and $\boldsymbol{\Psi}$ plays a pivotal role in determining the final rate of convergence. We also discuss specific structures of \mathbf{A} and $\boldsymbol{\Psi}$ to gain more insights. The results of this paper provide theoretical justifications for the popular practices in which sparse PCA is directly applied on high dimensional time series data for data visualization and feature selection. Examples areas include financial time series, biomedical imaging, and social media, etc.

The rest of the paper is organized as follows. In the next section, we briefly introduce sparse PCA and VAR time series model. In Section 3, we derive several useful results about the rate of convergence of sparse PCA. We prove the main theoretical result in Section 4. In section 5, we conduct numerical experiments on both simulated and real-world data to back up our theory.

2 Background

In this section, we briefly introduce the background of this paper. We start with notation. Let $\mathbf{A} = [\mathbf{A}_{jk}] \in \mathbb{R}^{d \times d}$ and $\mathbf{v} = (v_1, \dots, v_d)^T \in \mathbb{R}^d$. For $0 \leq q \leq 1$, we define the vector ℓ_q norms as $\|\mathbf{v}\|_q := (\sum_{j=1}^d |v_j|^q)^{\frac{1}{q}}$. Specifically $\|\mathbf{v}\|_0 = \text{card}\{\text{supp}(\mathbf{v})\}$. We denote $\|\mathbf{A}\|_q$ to be the operator norm of matrix \mathbf{A} . In particular, $\|\mathbf{A}\|_2$ is the spectral norm. Specifically, for $q = 1$ or $q = \infty$, $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq d} \sum_{i=1}^d |\mathbf{A}_{ij}|$, and $\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq d} \sum_{j=1}^d |\mathbf{A}_{ij}|$. We have $\|\mathbf{A}\|_2 \leq \sqrt{\|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty}$. Let $\lambda_j(\mathbf{A})$ be the j -th largest eigenvalue of \mathbf{A} . The d -dimension Euclidean unit sphere is $\mathbb{S}^{d-1} := \{\mathbf{v} | \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2 = 1\}$. The ℓ_q ball with radius R_q is $\mathbb{B}_q(R_q) := \{\mathbf{v} \in \mathbb{R}^d : \|\mathbf{v}\|_q \leq R_q\}$. For vector \mathbf{v}_1 and \mathbf{v}_2 , define inner product as $\langle \mathbf{v}_1, \mathbf{v}_2 \rangle := \mathbf{v}_1^T \mathbf{v}_2$. For matrix $\mathbf{A}_1, \mathbf{A}_2$, we define the inner product as $\langle \mathbf{A}_1, \mathbf{A}_2 \rangle := \text{tr}(\mathbf{A}_1^T \mathbf{A}_2)$. For a set \mathbb{K} , $|\mathbb{K}|$ is its cardinality. We use $\mathbf{0}$ to denote the all-zero matrix.

2.1 Vector Autoregressive Time Series

The weakly stationary Vector Autoregressive (VAR) time series model linear dependencies between different movements. In particular, the model assumes the T observations $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ are generated by the lag 1 autoregressive process:

$$\mathbf{X}_{t+1} = \mathbf{A}\mathbf{X}_t + \mathbf{Z}_t, \quad \text{for } t = 1, \dots, T-1. \quad (2.1)$$

To secure the weakly stationary of the above process, the transition matrix \mathbf{A} must have bounded spectral norm $\|\mathbf{A}\|_2 < 1$. We assume the Gaussian colored noise $\mathbf{Z}_1, \mathbf{Z}_2, \dots \stackrel{\text{i.i.d.}}{\sim} N_d(\mathbf{0}, \Psi)$. By assumption \mathbf{Z}_t and $\mathbf{X}_t, \mathbf{X}_{t-1}, \dots$ are independent. The stationary property indicates

$$\Sigma = \mathbf{A}\Sigma\mathbf{A}^T + \Psi. \quad (2.2)$$

For $j \geq k$, the covariance between \mathbf{X}_j and \mathbf{X}_k is

$$\text{Cov}(\mathbf{X}_j, \mathbf{X}_k) = \underbrace{\mathbf{A} \cdots \mathbf{A}}_{j-k} \Sigma := \mathbf{A}^{j-k} \Sigma.$$

In the sequel we call \mathbf{A} the transition matrix and Ψ the noise matrix. VAR model is widely used in the analysis of economic time series (Sims, 1980; Hatemi-J, 2004; Briiggemann and Lütkepohl, 2001), signal processing (de Waele and Broersen, 2003) and brain fMRI (Goebel et al., 2003; Roebroeck et al., 2005).

3 Sparse PCA for VAR Time Series

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ be the T observations of a random vector $\mathbf{X} \in \mathbb{R}^d$. Let $\mathbf{u}_1, \dots, \mathbf{u}_m$ be the top m eigenvectors of the covariance matrix Σ . Let $\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_m$ be the top m eigenvector of the sample covariance matrix \mathbf{S} . In low dimensions, PCA uses $\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_m$ to estimate $\mathbf{u}_1, \dots, \mathbf{u}_m$.

In the high dimensional settings where $d > T$, we assume that the leading eigenvector of Σ are under certain sparse constraints. In other words, we assume that \mathbf{u}_1 satisfies that $\|\mathbf{u}\|_0 \leq s$ and $\|\mathbf{u}\|_1 = 1$, i.e. $\mathbf{u}_1 \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(s)$. In this way, we define the model $\mathcal{M}(s, \Sigma, \lambda_1, \lambda_2)$ as follows:

$$\begin{aligned} \mathcal{M}(s, \Sigma, \lambda_1, \lambda_2) &:= \{\mathbf{X} : \mathbf{X} \sim N_d(\boldsymbol{\mu}, \Sigma), \\ &\lambda_1(\Sigma) = \lambda_1, \lambda_2(\Sigma) = \lambda_2, \|\mathbf{u}_1\|_0 = s\}, \end{aligned}$$

for some $\boldsymbol{\mu} \in \mathbb{R}^d$. We define $\hat{\mathbf{u}}_1$ to be the solution to the following optimization problem:

$$\begin{aligned} \hat{\mathbf{u}}_1 &:= \arg \max_{\mathbf{v} \in \mathbb{R}^d} \mathbf{v}^T \mathbf{S} \mathbf{v}, \\ &\text{subject to } \mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(s). \end{aligned} \quad (3.1)$$

Here $\hat{\mathbf{u}}_1$ is the global optimal estimator of \mathbf{u}_1 . In this paper we only discuss about $\hat{\mathbf{u}}_1$ corresponding to the leading eigenvalue $\lambda_1(\Sigma)$. Analysis for $\hat{\mathbf{u}}_k$ corresponding to $\lambda_k(\Sigma)$ is discussed in the longer version of this paper.

4 Theoretical Properties

In this section we provide the theoretical properties of the sparse PCA estimator for VAR time series. In particular, we provide the nonasymptotic upper bound of the rate of convergence in parameter estimation under the VAR model. To our knowledge, this is the first work analyzing the theoretical performance of PCA for the dependent data in high dimensions.

4.1 Main Result

The main result states that under the VAR model, the estimator $\hat{\mathbf{u}}_1$ obtained by (3.1) can approximate \mathbf{u}_1 in a parametric rate with respect to (n, d, s) , and the upper bound is also related to the transition matrix \mathbf{A} and the noisy matrix Ψ .

Theorem 4.1. *Provided that the random vector sequence $\{\mathbf{X}_t\}_{t=1}^T$ follows the VAR model described in (2.1) and $\mathbf{X}_t \in \mathcal{M}(s, \Sigma, \lambda_1, \lambda_2)$ for $t = 1, \dots, T$, the estimator $\hat{\mathbf{u}}_1$ derived in (3.1) has the following property:*

$$|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)| = O_P \left(\frac{1}{\lambda_1 - \lambda_2} \sqrt{\frac{s \log d}{T} \left(\frac{\|\Sigma\|_2}{1 - \|\mathbf{A}\|_2} \right)} \right).$$

Here for any two vectors \mathbf{v}_1 and $\mathbf{v}_2 \in \mathbb{S}^{d-1}$, $|\sin \angle(\mathbf{v}_1, \mathbf{v}_2)| := \sqrt{1 - (\mathbf{v}_1^T \mathbf{v}_2)^2}$.

Remark 4.2. *The bound obtained in (4.1) depends on Σ, \mathbf{A}, Ψ , where \mathbf{A} characterizes the data dependence degree. When both $\|\mathbf{A}\|_2$ and $\|\Psi\|_2$ do not scale with (n, d, s) , this is the parametric optimal rate (Ma, 2011; Vu and Lei, 2012).*

4.2 Technical Proofs

To prove Theorem 4.1, first we need to prove several lemmas. The following Lemma connects $\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)$ with $\sup_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(2s)} |\mathbf{v}^T (\Sigma - \mathbf{S}) \mathbf{v}|$. In the sequel, we assume that the assumptions in Theorem 4.1 hold.

Lemma 4.3. *\mathbf{u}_1 and $\hat{\mathbf{u}}_1$ satisfy*

$$\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1) \leq \frac{2}{\lambda_1 - \lambda_2} \sup_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(2s)} |\mathbf{v}^T (\Sigma - \mathbf{S}) \mathbf{v}|. \quad (4.1)$$

Proof. Let $\lambda_1 \geq \dots \geq \lambda_d$ be the eigenvalues of Σ . Let $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ be the corresponding eigenvectors. We have $\mathbf{u}_i^T \mathbf{u}_j = 0$ for $i \neq j$ and $\Sigma = \sum_{j=1}^d \lambda_j \mathbf{u}_j \mathbf{u}_j^T$. Let $\Sigma = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \Phi_0$, where Φ_0 can be represented by \mathbf{u}_1 and Σ as

$$\begin{aligned} \Phi_0 &= \Sigma - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T \\ &= \Sigma - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T \\ &= \Sigma - \mathbf{u}_1 \mathbf{u}_1^T \Sigma - \Sigma \mathbf{u}_1 \mathbf{u}_1^T + \mathbf{u}_1 (\mathbf{u}_1^T \Sigma \mathbf{u}_1) \mathbf{u}_1^T \\ &= (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \Sigma (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T). \end{aligned}$$

For any $\mathbf{u} \in \mathbb{S}^{d-1}$, we have

$$\begin{aligned} \langle \Sigma, \mathbf{u}_1 \mathbf{u}_1^T - \mathbf{u} \mathbf{u}^T \rangle &= \langle \Sigma, \mathbf{u}_1 \mathbf{u}_1^T \rangle - \langle \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \Phi_0, \mathbf{u} \mathbf{u}^T \rangle \\ &= \lambda_1 - \lambda_1 \langle \mathbf{u}_1, \mathbf{u} \rangle^2 - \langle \Phi_0, \mathbf{u} \mathbf{u}^T \rangle \\ &= \lambda_1 - \lambda_1 \langle \mathbf{u}_1, \mathbf{u} \rangle^2 - \mathbf{u}^T (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \Sigma (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{u}. \end{aligned}$$

Now we consider the unit vector

$$\mathbf{a} = (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{u} / \|(\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{u}\|_2 \in \mathbb{R}^d.$$

It is easy to verify \mathbf{a} is orthogonal to \mathbf{u}_1 . Therefore, $\mathbf{a} \in \text{span}\{\mathbf{u}_2, \dots, \mathbf{u}_d\}$. Letting $a_j = \mathbf{a}^T \mathbf{u}_j$, $j = 2, \dots, d$, we can get $\sum_{j=2}^d a_j^2 = 1$. Therefore we have

$$\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a} = \mathbf{a}^T \left(\lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \sum_{j=2}^d \lambda_j \mathbf{u}_j \mathbf{u}_j^T \right) \mathbf{a} = \sum_{j=2}^d \lambda_j a_j^2 \leq \lambda_2,$$

which indicates

$$\mathbf{u}^T (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \boldsymbol{\Sigma} (\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{u} \leq \lambda_2 \|(\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{u}\|_2^2.$$

Since \mathbf{u} is on the unit sphere \mathbb{S}^{d-1} , we have $\|\mathbf{u}\|_2 = 1$. Therefore, $\|(\mathbf{I}_d - \mathbf{u}_1 \mathbf{u}_1^T) \mathbf{u}\|_2^2 = 1 - \langle \mathbf{u}_1, \mathbf{u} \rangle^2$. Now we obtain

$$\langle \boldsymbol{\Sigma}, \mathbf{u}_1 \mathbf{u}_1^T - \mathbf{u} \mathbf{u}^T \rangle \geq (\lambda_1 - \lambda_2) (1 - \langle \mathbf{u}_1, \mathbf{u} \rangle^2). \quad (4.2)$$

Since (4.2) holds for any $\mathbf{u} \in \mathbb{S}^{d-1}$, letting $\mathbf{u} = \hat{\mathbf{u}}_1$ we have

$$\langle \boldsymbol{\Sigma}, \mathbf{u}_1 \mathbf{u}_1^T - \mathbf{u} \mathbf{u}^T \rangle \geq (\lambda_1 - \lambda_2) \sin^2 \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1).$$

Since $\hat{\mathbf{u}}_1$ is defined as $\hat{\mathbf{u}}_1 := \arg \max_{\mathbf{v} \in \mathbb{R}^d} \mathbf{v}^T \mathbf{S} \mathbf{v}$, we know

$$\hat{\mathbf{u}}_1^T \mathbf{S} \hat{\mathbf{u}}_1 - \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \leq 0.$$

Therefore, $\langle \mathbf{S}, \mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T \rangle \leq 0$, we have

$$\begin{aligned} \sin^2 \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1) &\leq \frac{1}{\lambda_1 - \lambda_2} \langle \boldsymbol{\Sigma}, \mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T \rangle \\ &\leq \frac{1}{\lambda_1 - \lambda_2} \langle \boldsymbol{\Sigma} - \mathbf{S}, \mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T \rangle. \end{aligned}$$

Let $\boldsymbol{\Pi}$ be the diagonal matrix with diagonal values being 1 if and only if the corresponding entries in \mathbf{u}_1 or $\hat{\mathbf{u}}_1$ are zero. Therefore, we know there are at most $2s$ nonzero values in $\boldsymbol{\Pi}$. Then we have

$$\begin{aligned} &\sin^2 \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1) \\ &\leq \frac{1}{\lambda_1 - \lambda_2} \langle \boldsymbol{\Sigma} - \mathbf{S}, \mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T \rangle \\ &= \frac{1}{\lambda_1 - \lambda_2} \langle \boldsymbol{\Sigma} - \mathbf{S}, \boldsymbol{\Pi} (\mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T) \boldsymbol{\Pi} \rangle \\ &= \frac{1}{\lambda_1 - \lambda_2} \langle \boldsymbol{\Pi} (\boldsymbol{\Sigma} - \mathbf{S}) \boldsymbol{\Pi}, \mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T \rangle \\ &\leq \frac{1}{\lambda_1 - \lambda_2} \|\boldsymbol{\Pi} (\boldsymbol{\Sigma} - \mathbf{S}) \boldsymbol{\Pi}\|_2 \|\mathbf{u}_1 \mathbf{u}_1^T - \hat{\mathbf{u}}_1 \hat{\mathbf{u}}_1^T\|_S \\ &= \frac{1}{\lambda_1 - \lambda_2} \|\boldsymbol{\Pi} (\boldsymbol{\Sigma} - \mathbf{S}) \boldsymbol{\Pi}\|_2 \cdot 2 |\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)|. \end{aligned}$$

Here $\|\cdot\|_S$ denotes the sum of singular values. This implies

$$\begin{aligned} |\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)| &\leq \frac{2}{\lambda_1 - \lambda_2} \|\boldsymbol{\Pi} (\boldsymbol{\Sigma} - \mathbf{S}) \boldsymbol{\Pi}\|_2 \\ &\leq \frac{2}{\lambda_1 - \lambda_2} \sup_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(2s)} |\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}|. \end{aligned}$$

This completes the proof. \square

The next lemma comes from Ledoux and Talagrand (2011) and is informative in the proof of the main theorem.

Lemma 4.4. *Provided that $x_1, x_2, \dots, x_T \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$. $\mathbf{X} = (x_1, x_2, \dots, x_T)^T \in \mathbb{R}^T$ is a random vector. Mapping $f : \mathbb{R}^T \rightarrow \mathbb{R}$ is Lipschitz, i.e., for any $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^T$:*

$$\exists L \geq 0, \quad \text{s.t. } |f(\mathbf{v}_1) - f(\mathbf{v}_2)| \leq L \|\mathbf{v}_1 - \mathbf{v}_2\|_2,$$

Then for any $t > 0$ we have,

$$\mathbb{P}\left(|f(\mathbf{X}) - \mathbb{E}f(\mathbf{X})| > t\right) \leq 2 \exp\left(-\frac{t^2}{2L^2}\right). \quad (4.3)$$

Proof. We refer to Ledoux and Talagrand (2011)'s proof of this lemma. \square

The next lemma quantifies the difference between $\boldsymbol{\Sigma}$ and \mathbf{S} for any fixed vector $\mathbf{v} \in \mathbb{R}^d$.

Lemma 4.5. *Letting $\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_q(R_q)$, we have*

$$\left| \sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} - \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}} \right| = O_P\left(\sqrt{\frac{1}{T}}\right). \quad (4.4)$$

Proof. Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ be the T observations. $\mathbf{Y} = (\mathbf{x}_1^T \mathbf{v}, \mathbf{x}_2^T \mathbf{v}, \dots, \mathbf{x}_T^T \mathbf{v})^T \in \mathbb{R}^T$ is a zero-mean Gaussian random vector. We denote the covariance matrix of \mathbf{Y} as $\boldsymbol{\Sigma}_{\mathbf{Y}}$. $\boldsymbol{\Sigma}_{\mathbf{Y}}$ is symmetric and semi-definite. Thus $\boldsymbol{\Sigma}_{\mathbf{Y}}$ can be decomposed as $\boldsymbol{\Sigma}_{\mathbf{Y}} = \mathbf{Q}^T \mathbf{Q}$, where \mathbf{Q} is a matrix with orthogonal columns. Let $\sigma = \|\mathbf{Q}\|_2 = \sqrt{\|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_2}$. According to the definition of \mathbf{Y} , we have

$$\mathbf{v}^T \mathbf{S} \mathbf{v} = \mathbf{v}^T \left(\frac{\sum_{i=1}^T \mathbf{x}_i \mathbf{x}_i^T}{T} \right) \mathbf{v} = \frac{\mathbf{Y}^T \mathbf{Y}}{T} = \left(\frac{\|\mathbf{Y}\|_2}{\sqrt{T}} \right)^2,$$

$$\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} = \mathbf{v}^T \mathbb{E}(\mathbf{X} \mathbf{X}^T) \mathbf{v} = \mathbb{E} \left(\frac{\mathbf{Y}^T \mathbf{Y}}{T} \right) = \mathbb{E} \left(\frac{\|\mathbf{Y}\|_2}{\sqrt{T}} \right).$$

For convenience we define $\mathbf{W} := \|\mathbf{Y}\|_2 / \sqrt{T}$ and $f(\mathbf{v}) := \|\mathbf{Q} \mathbf{v}\|_2 / \sqrt{T}$ for $\mathbf{v} \in \mathbb{R}^T$. Since we have $\mathbf{Y} = \mathbf{Q} \mathbf{V}$, $\mathbf{V} \in \mathbb{R}^T$ is a zero-mean Gaussian vector with covariance matrix \mathbf{I}_T . It can be verified that mapping $f : \mathbb{R}^T \rightarrow \mathbb{R}$ is Lipschitz with $L = \sigma / \sqrt{T}$. Using Lemma 4.4, we can get

$$\mathbb{P}\left(\|\mathbf{W} - \mathbb{E}(\mathbf{W})\| \geq t\right) \leq 2 \exp\left(-\frac{t^2 T}{2\sigma^2}\right). \quad (4.5)$$

Since $\text{Var}(\mathbf{W}) = \mathbb{E}(\mathbf{W}^2) - \mathbb{E}^2(\mathbf{W}) \geq 0$, at the same time we have $\mathbb{E}(\mathbf{W}^2) \geq 0, \mathbb{E}(\mathbf{W}) \geq 0$, which implies $\sqrt{\mathbb{E}(\mathbf{W}^2)} - \mathbb{E}(\mathbf{W}) \geq 0$. Thus we get

$$\left(\sqrt{\mathbb{E}(\mathbf{W}^2)} - \mathbb{E}(\mathbf{W}) \right)^2 \leq \mathbb{E}\left(\|\mathbf{W} - \mathbb{E}(\mathbf{W})\|^2\right) = \frac{4\sigma^2}{T},$$

which indicates

$$\left| \sqrt{\mathbb{E}(\mathbf{W}^2)} - \mathbb{E}(\mathbf{W}) \right| \leq \frac{2\sigma}{\sqrt{T}}. \quad (4.6)$$

Here (4.6) together with $\|\mathbf{W} - \mathbb{E}(\mathbf{W})\| \leq t$ implies $\left| \mathbf{W} - \sqrt{\mathbb{E}(\mathbf{W}^2)} \right| \leq t + 2\sigma / \sqrt{T}$. Therefore, according

to the definition of \mathbf{Y} and \mathbf{W} ,

$$\begin{aligned} & \mathbb{P}\left(\left|\sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} - \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}}\right| \geq t + \frac{2\sigma}{\sqrt{T}}\right) \\ &= \mathbb{P}\left(\left|\mathbf{W} - \sqrt{\mathbb{E}(\mathbf{W}^2)}\right| \geq t + \frac{2\sigma}{\sqrt{T}}\right) \\ &\leq \mathbb{P}\left(\|\mathbf{W} - \mathbb{E}(\mathbf{W})\| \geq t\right) \\ &\leq 2 \exp\left(-\frac{t^2 T}{2\sigma^2}\right). \end{aligned} \quad (4.7)$$

We reach the conclusion

$$\left|\sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} - \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}}\right| = O_p\left(\sqrt{\frac{1}{T}}\right).$$

This completes the proof. \square

Lemma 4.6. *An ϵ -net \mathbb{N}_ϵ of a sphere \mathbb{S}^{n-1} (equipped with Euclidean distance) is a subset of \mathbb{S}^{n-1} such that for any $\mathbf{v} \in \mathbb{S}^{n-1}$, there exists $\mathbf{u} \in \mathbb{N}_\epsilon$ subject to $\|\mathbf{u} - \mathbf{v}\|_2 \leq \epsilon$. It is shown by Vershynin (2010) that for any $\epsilon > 0$,*

$$|\mathbb{N}_\epsilon| \leq \left(1 + \frac{2}{\epsilon}\right)^n. \quad (4.8)$$

Also, for matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, the inequality below holds for any $\epsilon \in [0, 1]$

$$\max_{\mathbf{v}_1 \in \mathbb{S}^{n-1}} |\mathbf{v}_1 \mathbf{A} \mathbf{v}_1| \leq (1 - 2\epsilon)^{-1} \max_{\mathbf{v}_2 \in \mathbb{N}_\epsilon} |\mathbf{v}_2 \mathbf{A} \mathbf{v}_2|. \quad (4.9)$$

Proof. Construct a maximal ϵ -separated subset \mathbb{N}_ϵ of \mathbb{S}^{n-1} . Here an ϵ -separated set is defined as the set whose arbitrary two different elements \mathbf{x}, \mathbf{y} are at least distance ϵ away, i.e. $\|\mathbf{x} - \mathbf{y}\|_2 \geq \epsilon$. Since \mathbb{N}_ϵ is a maximal ϵ -separated subset of \mathbb{S}^{n-1} , there isn't any other ϵ -separated \mathbb{N}'_ϵ of \mathbb{S}^{n-1} such that $\mathbb{N}_\epsilon \subset \mathbb{N}'_\epsilon$.

We can prove that \mathbb{N}_ϵ is an ϵ -net of \mathbb{S}^{n-1} by contradiction. If we assume for $\mathbf{u} \in \mathbb{S}^{n-1}$, there is a point $\mathbf{v} \in \mathbb{N}_\epsilon$, $\|\mathbf{u} - \mathbf{v}\|_2 > \epsilon$, then $\mathbb{N}'_\epsilon = \{\mathbf{v}\} \cup \mathbb{N}_\epsilon$ is a larger ϵ -separated subset of \mathbb{S}^{n-1} that contains \mathbb{N}_ϵ , which contradicts with the fact that \mathbb{N}_ϵ is the maximal ϵ -separated subset of \mathbb{S}^{n-1} .

Now we derive the bound of $|\mathbb{N}_\epsilon|$. We cover the neighborhood of every $\mathbf{x}_i \in \mathbb{N}_\epsilon$ with disjoint balls $\mathbb{B}_{\mathbf{y}_i} = \{\mathbf{y}_i | \|\mathbf{y}_i - \mathbf{x}_i\|_2 < \epsilon/2\}$. For any $\mathbb{B}_{\mathbf{y}_i} \subset \mathbb{B}_0 = \{\mathbf{y} | \|\mathbf{y}\|_2 < 1 + \epsilon/2\}$ we have

$$\begin{aligned} |\mathbb{N}_\epsilon| \cdot |\mathbb{B}_{\mathbf{y}_i}| &= \sum_{i=1}^{|\mathbb{N}_\epsilon|} |\mathbb{B}_{\mathbf{y}_i}| = \left| \bigcup_{i=1}^{|\mathbb{N}_\epsilon|} \mathbb{B}_{\mathbf{y}_i} \right| \leq |\mathbb{B}_0|, \\ |\mathbb{N}_\epsilon| &\leq \frac{|\mathbb{B}_0|}{|\mathbb{B}_{\mathbf{y}_i}|} = \frac{(1 + \frac{\epsilon}{2})^n}{(\frac{\epsilon}{2})^n} = \left(1 + \frac{2}{\epsilon}\right)^n. \end{aligned}$$

Then we turn to prove (4.9). For $\mathbf{v}_1 \in \mathbb{S}^{n-1}$ and $\mathbf{v}_2 \in \mathbb{N}_\epsilon$, we know since $\|\mathbf{v}_1 - \mathbf{v}_2\|_2 \leq \epsilon$,

$$\begin{aligned} |\mathbf{v}_1 \mathbf{A} \mathbf{v}_1 - \mathbf{v}_2 \mathbf{A} \mathbf{v}_2| &\leq \|\mathbf{A}\|_2 \|\mathbf{v}_1\|_2 \|\mathbf{v}_1 - \mathbf{v}_2\|_2 \\ &\quad + \|\mathbf{A}\|_2 \|\mathbf{v}_2\|_2 \|\mathbf{v}_1 - \mathbf{v}_2\|_2 \leq 2\epsilon \|\mathbf{A}\|_2 \end{aligned}$$

It follows that

$$|\mathbf{v}_2 \mathbf{A} \mathbf{v}_2| \geq (1 - 2\epsilon) \|\mathbf{A}\|_2 = (1 - 2\epsilon) \max_{\mathbf{v}_1 \in \mathbb{S}^{d-1}} |\mathbf{v}_1 \mathbf{A} \mathbf{v}_1|$$

Taking the maximum over $\mathbf{v}_2 \in \mathbb{N}_\epsilon$, we completes the proof. \square

Combining the above Lemmas, we can now proceed to the main proof.

Proof of Theorem 4.1. Using Lemma 4.5, it is easy to show

$$\sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} = \sqrt{\mathbb{E}(\mathbf{W}^2)} = \sqrt{\frac{\text{tr}(\mathbf{Q}\mathbf{Q}^T)}{T}} = \frac{\|\mathbf{Q}\|_F}{\sqrt{T}} \leq \sigma.$$

Then using (4.7) we can get

$$\begin{aligned} & \mathbb{P}\left(\left|\sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} + \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}}\right| \geq t + 4\sigma\right) \\ &\leq \mathbb{P}\left(\left|\sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}} + \sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}}\right| \geq t + \frac{2\sigma}{\sqrt{T}} + 2\sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}}\right) \\ &\leq \mathbb{P}\left(\left|\sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} - \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}}\right| \geq t + \frac{2\sigma}{\sqrt{T}}\right) \\ &\leq 2 \exp\left(-\frac{t^2 T}{2\sigma^2}\right). \end{aligned} \quad (4.10)$$

We define events \mathcal{E}_1 and \mathcal{E}_2 , combining (4.7) (with $t = t_1$) and (4.10) (with $t = t_2$),

$$\begin{aligned} \mathcal{E}_1 &:= \left\{ \left| \sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} - \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}} \right| \geq t_1 + 2\frac{\sigma}{\sqrt{T}} \right\} \\ \mathcal{E}_2 &:= \left\{ \left| \sqrt{\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}} + \sqrt{\mathbf{v}^T \mathbf{S} \mathbf{v}} \right| \geq t_2 + 4\sigma \right\}. \end{aligned}$$

We derive for any $t_1, t_2 > 0$,

$$\begin{aligned} & \mathbb{P}\left(\left|\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}\right| \geq \left(t_1 + 2\frac{\sigma}{\sqrt{T}}\right) (t_2 + 4\sigma)\right) \\ &\leq \mathbb{P}(\mathcal{E}_1) + \mathbb{P}(\mathcal{E}_2) \\ &\leq 2 \exp\left(-\frac{t_1^2 T}{2\sigma^2}\right) + 2 \exp\left(-\frac{t_2^2 T}{2\sigma^2}\right). \end{aligned} \quad (4.11)$$

Now we turn to upper bound $\sup_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(2\sigma)} |\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}|$ in Lemma 4.3. Assuming we have a fixed subset $\mathbb{K} \subset \{1, \dots, d\}$, we define

$$\mathbb{B}_{\mathbb{K}} = \{\mathbf{v} | \text{for any } i \in \{1, \dots, d\} \setminus \mathbb{K}, v_i = 0\}.$$

For any $t_1, t_2 > 0$, we define event $\mathcal{E}_{\mathbb{K}}$ and $\mathcal{E}_{\mathbf{v}}$ as

$$\begin{aligned} \mathcal{E}_{\mathbb{K}} &:= \left\{ \max_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_{\mathbb{K}}} |\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}| \geq 2 \left(t_1 + 2\frac{\sigma}{\sqrt{T}}\right) (t_2 + 4\sigma) \right\}, \\ \mathcal{E}_{\mathbf{v}} &:= \left\{ \left|\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}\right| \geq \left(t_1 + 2\frac{\sigma}{\sqrt{T}}\right) (t_2 + 4\sigma) \right\}. \end{aligned}$$

According to Lemma 4.6, we define the $\frac{1}{4}$ -net of $\mathbb{S}^{d-1} \cap \mathbb{B}_{\mathbb{K}}$ as $\mathbb{N}_{\mathbb{K}}$, then we can get

$$\mathcal{E}_{\mathbb{K}} \subset \bigcup_{\mathbf{v} \in \mathbb{N}_{\mathbb{K}}} \left\{ \left|\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}\right| \geq \left(t_1 + 2\frac{\sigma}{\sqrt{T}}\right) (t_2 + 4\sigma) \right\},$$

Combining (4.11) and (4.8) and letting $|\mathbb{K}| = 2s$, we obtain

$$\begin{aligned} \mathbb{P}(\mathcal{E}_{\mathbb{K}}) &\leq \sum_{\mathbf{v} \in \mathbb{N}_{\mathbb{K}}} \mathbb{P}(\mathcal{E}_{\mathbf{v}}) = |\mathbb{N}_{\mathbb{K}}| \mathbb{P}(\mathcal{E}_{\mathbf{v}}) \\ &\leq 9^{2s} \left(2 \exp\left(-\frac{t_1^2 T}{2\sigma^2}\right) + 2 \exp\left(-\frac{t_2^2 T}{2\sigma^2}\right) \right). \end{aligned}$$

Now we consider arbitrary subset $\mathbb{K} \subset \{1, \dots, d\}$ with cardinality $2s$. We define

$$\mathcal{E}'_{\mathbb{K}} := \left\{ \max_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(2s)} |\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}| \geq 2 \left(t_1 + 2 \frac{\sigma}{\sqrt{T}} \right) (t_2 + 4\sigma) \right\}.$$

Then we have

$$\begin{aligned} \mathbb{P}(\mathcal{E}'_{\mathbb{K}}) &\leq \sum_{\mathbb{K} \subset \{1, \dots, d\}} \mathbb{P}(\mathcal{E}_{\mathbb{K}}) \leq \binom{d}{2s} \mathbb{P}(\mathcal{E}_{\mathbb{K}}) \\ &\leq 9^{2s} \binom{d}{2s} \left(2 \exp\left(-\frac{t_1^2 T}{2\sigma^2}\right) + 2 \exp\left(-\frac{t_2^2 T}{2\sigma^2}\right) \right). \end{aligned}$$

Using (4.1), we can get for any $t_1, t_2 > 0$, $\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(2s)$, which implies

$$\begin{aligned} \mathbb{P}\left(|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)| \geq \frac{4}{\lambda_1 - \lambda_2} \left(t_1 + 2 \frac{\sigma}{\sqrt{T}} \right) (t_2 + 4\sigma) \right) \\ \leq \mathbb{P}\left(\max |\mathbf{v}^T (\boldsymbol{\Sigma} - \mathbf{S}) \mathbf{v}| \geq 2 \left(t_1 + 2 \frac{\sigma}{\sqrt{T}} \right) (t_2 + 4\sigma) \right) \\ \leq 9^{2s} \binom{d}{2s} \left(2 \exp\left(-\frac{t_1^2 T}{2\sigma^2}\right) + 2 \exp\left(-\frac{t_2^2 T}{2\sigma^2}\right) \right). \end{aligned}$$

Now we are going to derive the upper bound for $\sigma^2 = \|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_2$. We note that VAR model assumes $\|\mathbf{A}\|_2 < 1$. Thus σ^2 can be bounded as follows,

$$\sigma^2 = \|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_2 \leq \|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_{\infty} = \max_{i \in \{1, \dots, T\}} \sum_{j=1}^T |(\boldsymbol{\Sigma}_{\mathbf{Y}})_{ij}|. \quad (4.12)$$

We have $\|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_2 \leq \sqrt{\|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_1 \|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_{\infty}} = \|\boldsymbol{\Sigma}_{\mathbf{Y}}\|_{\infty}$, since $\boldsymbol{\Sigma}_{\mathbf{Y}}$ is symmetric. According to the definition (1.1) we have,

$$\mathbf{v}^T \text{Cov}(\mathbf{x}_i, \mathbf{x}_j) \mathbf{v} = \mathbf{v}^T \text{Cov}(\mathbf{x}_j, \mathbf{x}_i) \mathbf{v} = \mathbf{v}^T \mathbf{A}^{|i-j|} \boldsymbol{\Sigma} \mathbf{v},$$

In order to get $\max_{i \in \{1, \dots, T\}} \sum_{j=1}^T |(\boldsymbol{\Sigma}_{\mathbf{Y}})_{ij}|$ in (4.12), we first derive the upper bound of $\sum_{j=1}^T |(\boldsymbol{\Sigma}_{\mathbf{Y}})_{ij}|$.

$$\begin{aligned} \sum_{j=1}^T |(\boldsymbol{\Sigma}_{\mathbf{Y}})_{ij}| &= \sum_{j \neq i} |\mathbf{v}^T \text{Cov}(\mathbf{x}_i, \mathbf{x}_j) \mathbf{v}| + |\mathbf{v}^T \text{Cov}(\mathbf{x}_i, \mathbf{x}_i) \mathbf{v}| \\ &= \sum_{j \neq i} |\mathbf{v}^T \mathbf{A}^{|i-j|} \boldsymbol{\Sigma} \mathbf{v}| + |\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}| \\ &\leq \sum_{j \neq i} \|\boldsymbol{\Sigma}\|_2 \|\mathbf{A}\|_2^{|i-j|} + \|\boldsymbol{\Sigma}\|_2 \\ &\leq \frac{2\|\boldsymbol{\Sigma}\|_2}{1 - \|\mathbf{A}\|_2}. \end{aligned}$$

Thus, letting $\delta^2 = 2\|\boldsymbol{\Sigma}\|_2 / (1 - \|\mathbf{A}\|_2) \geq \sigma^2$, we have

for any $t_1, t_2 > 0$,

$$\begin{aligned} \mathbb{P}\left(\left| \sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1) \right| \geq \frac{4}{\lambda_1 - \lambda_2} \left(t_1 + 2 \frac{\delta}{\sqrt{T}} \right) (t_2 + 4\delta) \right) \\ \leq 9^{2s} \binom{d}{2s} \left(2 \exp\left(-\frac{t_1^2 T}{2\delta^2}\right) + 2 \exp\left(-\frac{t_2^2 T}{2\delta^2}\right) \right). \end{aligned}$$

Letting $t_1 = \sqrt{s\delta^2 \log d/T}$ and t_2 be a constant, we reach the conclusion,

$$|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)| = O_P\left(\frac{1}{\lambda_1 - \lambda_2} \sqrt{\frac{s \log d}{T} \left(\frac{\|\boldsymbol{\Sigma}\|_2}{1 - \|\mathbf{A}\|_2} \right)} \right).$$

This is the main result of the rate of convergence. \square

5 Experiments

In this section we show some experimental results on both synthetic and real-world data to back up the theoretical results we obtain in last section. We use the truncated power method proposed by Yuan and Zhang (2011) to approximate the global estimator $\hat{\mathbf{u}}_1$ obtained in (4.2).

5.1 Synthetic Data

In this section we experiment with sparse PCA on synthetic data. We show how \mathbf{A} , $\boldsymbol{\Sigma}$ and $\boldsymbol{\Psi}$ affect $|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)|$. First we create \mathbf{A} given $\|\mathbf{A}\|_2$. With $\lambda_1(\mathbf{A})$ and $\lambda_2(\mathbf{A})$, we generate

$$\boldsymbol{\Sigma} = (\lambda_1(\mathbf{A}) - 1) \mathbf{u}_1^T \mathbf{u}_1 + (\lambda_2(\mathbf{A}) - 1) \mathbf{u}_2^T \mathbf{u}_2 + \mathbf{I},$$

where $\|\mathbf{u}_1\|_0 = s$, \mathbf{u}_1 and \mathbf{u}_2 is orthogonal. According to the stationary property (2.2), the covariance matrix of the noise random vector \mathbf{Z}_i is $\boldsymbol{\Psi} = \boldsymbol{\Sigma} - \mathbf{A}^T \boldsymbol{\Sigma} \mathbf{A}$, where $\boldsymbol{\Psi}$ must be a positive semidefinite matrix. We generate $T = 50$ data points according to the VAR time series model in (1.1). We illustrate how the scaling of $\|\mathbf{A}\|_2$ affects the accuracy of estimator $\hat{\mathbf{u}}_1$. Set $\lambda_1(\boldsymbol{\Sigma}) = 10$, $\lambda_2(\boldsymbol{\Sigma}) = 5$, $s = 20$, $d = 200$ and $\|\mathbf{A}\|_2 \in [0, 0.9]$, we repeatedly experiment for 3000 times for each $\|\mathbf{A}\|_2$. Here $\boldsymbol{\Sigma}$ is fixed while \mathbf{A} and $\boldsymbol{\Psi}$ are varying. The results are illustrated in Figure 1 by plotting the relevant part in the theoretical upper bound $\frac{\|\boldsymbol{\Sigma}\|_2}{1 - \|\mathbf{A}\|_2}$ against the empirical error $|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)|$. As can be observed in Fig. 1, the empirical error increases when the spectral norm of transition matrix $\|\mathbf{A}\|_2$ increases. This makes sense because when $\|\mathbf{A}\|_2 = 0$, VAR model is reduced to the i.i.d. case, since \mathbf{X}_{t+1} doesn't depend on \mathbf{X}_t . When $\|\mathbf{A}\|_2 \rightarrow 1$, the degree of dependency increases and sparse PCA loses its estimation accuracy, as quantified in Theorem 4.1. Table 1 shows the values and the standard deviations of $|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)|$ for each $\|\mathbf{A}\|_2$.

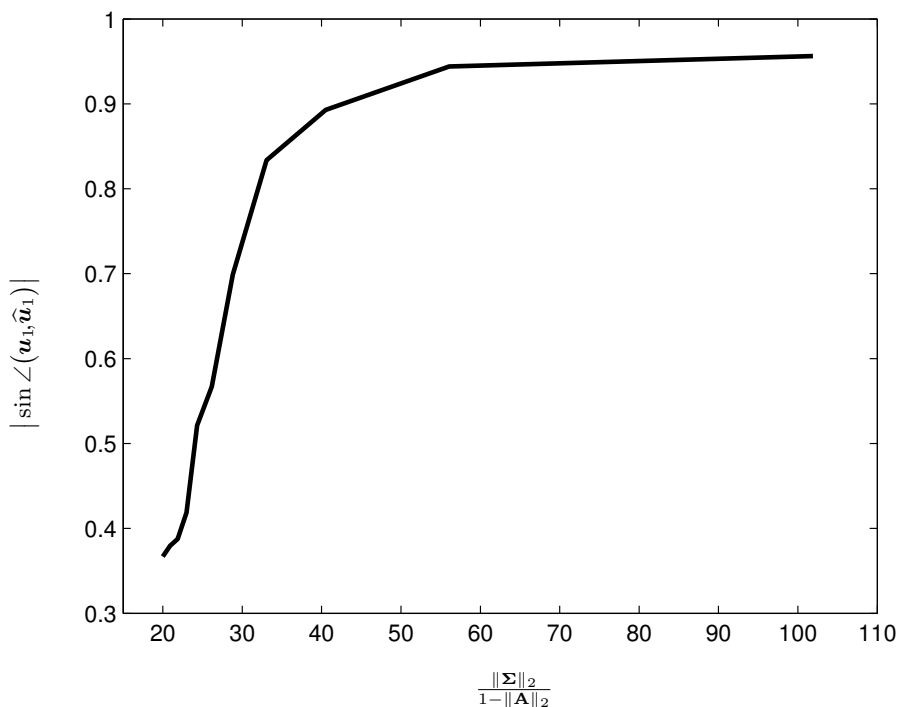


Figure 1: $\|\mathbf{A}\|_2$'s impact on the estimation accuracy $|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)|$. When dependency $\|\mathbf{A}\|_2$ increases toward 1, the upper bound $\frac{\|\Sigma\|_2}{1-\|\mathbf{A}\|_2}$ increases, thus the estimator $\hat{\mathbf{u}}_1$ ' accuracy decreases, i.e. $|\sin \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)|$ increases towards 1.

Table 1: Corresponding values when $\|\mathbf{A}\|_2$ increases. Here $\kappa^* = \frac{\|\Sigma\|_2}{1-\|\mathbf{A}\|_2}$ and $\theta = \angle(\mathbf{u}_1, \hat{\mathbf{u}}_1)$

$\ \mathbf{A}\ _2$	κ^*	$ \sin \theta $	$\sigma(\sin \theta)$
0	20.0000	0.3669	0.1284
0.0900	20.9080	0.3795	0.2150
0.1800	21.8711	0.3876	0.1319
0.2700	22.9696	0.4186	0.1591
0.3600	24.3290	0.5214	0.2201
0.4500	26.1568	0.5671	0.2625
0.5400	28.8231	0.6989	0.1981
0.6300	33.0580	0.8336	0.1549
0.7200	40.5303	0.8930	0.1410
0.8100	56.0706	0.9441	0.0765
0.9000	101.9000	0.9563	0.1196

5.2 Equity Data

In this section we apply sparse PCA on daily closing prices of 452 stocks in the S&P 500 index between January 1, 2003 through January 1, 2008 from Yahoo! Finance (finance.yahoo.com). That is to say, we have altogether $T = 1257$ observations corresponding to the vector of closing prices on a trading day. We categorize the 452 stocks into 10 Global Industry Classification Standard (GICS) sectors, including **Consumer Discretionary** (70 stocks), **Consumer Staples** (35 stocks), **Energy** (37 stocks), **Financials** (74 stocks), **Health Care** (46 stocks), **Industrials** (59 stocks), **Information Technology** (64 stocks), **Telecommunications Services** (6 stocks), **Materials** (29 stocks), and **Utilities** (32 stocks). We expect that the stocks from the same sector are likely to appear in the non-zero entries of the same principal component, since stocks from the same sector tend to be more correlated.

Let $\mathbf{P} = [\mathbf{P}_{t,j}]$ be the closing price of stock j on day t . In this paper we are interested in the transformed data, where we calculate the log-ratio of the price at time t to price at $t - 1$:

$$\mathbf{X}_{t,j} = \log(\mathbf{P}_{t+1,j}/\mathbf{P}_{t,j}), \quad t = 1, \dots, T - 1.$$

It is obvious that there exists data dependence struc-

ture between $\mathbf{X}_{t_1,j}$ and $\mathbf{X}_{t_2,j}$ for any $t_1, t_2 \leq T - 1$ and it accordingly raises concern for conducting classical sparse PCA algorithms on this dataset \mathbf{X} . However, the argument in this paper provides justifications to such a procedure. In particular, we conclude that the same parametric rate can be persisted if the operator norm of the transition matrix does not scale with (n, d, s) . Here we present labels of the first three estimated leading eigenvectors in Table 2. As can be observed, the sparse PCA algorithm tend to group stocks from the same sector into the same eigenvector.

Acknowledgements

This research was supported by NSF award IIS-1116730.

Table 2: Non-zero terms' sectors in the 1st, 2nd and 3rd eigenvectors obtained.

$\hat{\mathbf{u}}_1$	$\hat{\mathbf{u}}_2$	$\hat{\mathbf{u}}_3$
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Consumer Discretionary
Financials	Industrials	Financials
Financials	Industrials	Financials
Financials	Industrials	Financials
Financials	Industrials	Financials
Financials	Industrials	Financials
Financials	Industrials	Financials
Financials	Industrials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Industrials
Financials	Materials	Information Technology
Financials	Materials	Information Technology
Financials	Materials	Materials
Financials	Materials	Materials
Financials	Materials	Materials

References

- BRIIGEMANN, R. and LIITKEPOHL, H. (2001). Lag selection in subset var models with an application to a us monetary system. *Econometric Studies: A Festschrift in Honour of Joachim Frohn*, LIT-Verlag, Münster 107–28.
- D’ASPREMONT, A., BACH, F. and GHAOUI, L. (2008). Optimal solutions for sparse principal component analysis. *The Journal of Machine Learning Research* **9** 1269–1294.
- D’ASPREMONT, A., EL GHAOUI, L., JORDAN, M. and LANCKRIET, G. (2004). *A direct formulation for sparse PCA using semidefinite programming*. Computer Science Division, University of California.
- D’ASPREMONT, A., EL GHAOUI, L., JORDAN, M. and LANCKRIET, G. (2005). Sparse pca with applications in finance .
- DE WAELE, S. and BROERSEN, P. (2003). Order selection for vector autoregressive models. *Signal Processing, IEEE Transactions on* **51** 427–433.
- FAN, J., QI, L. and TONG, X. (2012). Penalized least squares estimation with weakly dependent data .
- GOEBEL, R., ROEBROECK, A., KIM, D. and FORMISANO, E. (2003). Investigating directed cortical interactions in time-resolved fmri data using vector autoregressive modeling and granger causality mapping. *Magnetic resonance imaging* **21** 1251–1261.
- HATEMI-J, A. (2004). Multivariate tests for autocorrelation in the stable and unstable var models. *Economic Modelling* **21** 661–683.
- JOHNSTONE, I. and LU, A. (2009). On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association* **104** 682–693.
- JOLLIFFE, I., TRENDAFILOV, N. and UDDIN, M. (2003). A modified principal component technique based on the lasso. *Journal of Computational and Graphical Statistics* **12** 531–547.
- JOURNÉE, M., NESTEROV, Y., RICHTÁRIK, P. and SEPULCHRE, R. (2010). Generalized power method for sparse principal component analysis. *The Journal of Machine Learning Research* **11** 517–553.
- LEDoux, M. and TALAGRAND, M. (2011). *Probability in Banach Spaces: isoperimetry and processes*, vol. 23. Springer.
- LOH, P. and WAINWRIGHT, M. (2011). High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. *Proceedings of the Twenty-Third Annual Conference on Neural Information Processing Systems (NIPS)* .
- MA, Z. (2011). Sparse principal component analysis and iterative thresholding. *Arxiv preprint arXiv:1112.2432* .
- ROEBROECK, A., FORMISANO, E., GOEBEL, R. ET AL. (2005). Mapping directed influence over the brain using granger causality and fmri. *Neuroimage* **25** 230–242.
- SHEN, H. and HUANG, J. (2008). Sparse principal component analysis via regularized low rank matrix approximation. *Journal of multivariate analysis* **99** 1015–1034.
- SIMS, C. A. (1980). Macroeconomics and reality. *Econometrica* **48** 1–48.
URL <http://ideas.repec.org/a/econ/emetrp/v48y1980i1p1-48.html>
- SKINNER, C., HOLMES, D. and SMITH, T. (1986). The effect of sample design on principal component analysis. *Journal of the American Statistical Association* **81** 789–798.
- VERSHYNIN, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027* .
- VU, V. and LEI, J. (2012). Minimax rates of estimation for sparse pca in high dimensions. *Arxiv preprint arXiv:1202.0786* .
- WITTEN, D., TIBSHIRANI, R. and HASTIE, T. (2009). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* **10** 515–534.
- YUAN, M. (2010). High dimensional inverse covariance matrix estimation via linear programming. *The Journal of Machine Learning Research* **99** 2261–2286.
- YUAN, X. and ZHANG, T. (2011). Truncated power method for sparse eigenvalue problems. *Arxiv preprint arXiv:1112.2679* .
- ZHANG, Y., D’ASPREMONT, A. and GHAOUI, L. (2012). Sparse pca: Convex relaxations, algorithms and applications. *Handbook on Semidefinite, Conic and Polynomial Optimization* 915–940.
- ZHANG, Y. and EL GHAOUI, L. (2011). Large-scale sparse principal component analysis with application to text data. In *Advances in Neural Information Processing Systems (NIPS)*.
- ZOU, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* **101** 1418–1429.