

An analysis of the SPARSEVA estimate for the finite sample data case [★]

Huong Ha ^a, James S. Welsh ^a, Cristian R. Rojas ^b, Bo Wahlberg ^b

^a*School of Electrical Engineering and Computer Science, The University of Newcastle, Australia*

^b*Department of Automatic Control and ACCESS, School of Electrical Engineering, KTH Royal Institute of Technology, SE-100 44 Stockholm, Sweden*

Abstract

In this paper, we develop an upper bound for the SPARSEVA (SPARSe Estimation based on a VALidation criterion) estimation error in a general scheme, i.e., when the cost function is strongly convex and the regularized norm is decomposable for a pair of subspaces. We show how this general bound can be applied to a sparse regression problem to obtain an upper bound for the traditional SPARSEVA problem. Numerical results are used to illustrate the effectiveness of the suggested bound.

Key words: SPARSEVA estimate; upper bound; finite sample data.

1 Introduction

Regularization is a well known technique for estimating model parameters from measured input-output data. Its applications are in any fields that are related to constructing mathematical models from observed data, such as system identification, machine learning and econometrics. The idea of the regularization technique is to solve a convex optimization problem constructed from a cost function and a weighted regularizer (regularized M-estimators). There are various types of regularizers that have been suggested so far, such as the l_1 [19], l_2 [20] and nuclear norms [5] [6].

During the last few decades, in the system identification community, regularization has been utilised extensively [15], to impose properties of smoothness and sparsity in the estimated models (see, e.g., [13, 22]). Most of this work has focused on analysing the asymptotic properties of an estimator, i.e., when the length of the data goes to infinity. The purpose of this type of analysis is to evaluate the performance of the estimation method to determine if the estimate is acceptable. However, in practice, the data sample size for any estimation problem is always finite, hence, it is difficult to judge the performance of the estimated parameters based on asymptotic properties, especially when the data length is short.

Recently, a number of authors have published research ([1], [4], [12]) aimed at analysing estimation error properties of the regularized M-estimators when the sample size of the data is finite. Specifically, they develop upper bounds on the estimation error for high dimensional problems, i.e., when the number of parameters is comparable to or larger than the sample size of the data. Most of these activities are from the statistics and machine learning communities. Among these works, the paper [12] provides a very elegant and interesting framework for establishing consistency and convergence rates of estimates obtained from a regularized procedure under high dimensional scaling. It determines a general upper bound for regularized M-estimators and then shows how it can be used to derive bounds for some specific scenarios.

Here in this paper we utilize the framework suggested in [12] to develop an upper bound for the estimation error of the M-estimators used in a system identification problem. Here, the M-estimator problems are implemented using the SPARSEVA (SPARSe Estimation based on a VALidation criterion) framework [16], [17]. The approach in [12] has been developed for penalized estimators, so it has to be suitably modified for SPARSEVA, which is not a penalized estimator, but the solution of a constrained optimization problem. Our aim is to derive an upper bound for the estimation error of the general SPARSEVA estimate. We then apply this bound to a sparse linear regression problem to obtain an upper bound for the traditional SPARSEVA problem with some assumptions on the regression matrix. These assumptions can be considered as the price in order to derive the upper bound. In addition, we also provide numerical simulation results to illustrate the suggested bound of the SPARSEVA estimation

[★] The material in this paper was not presented at any conference.

Email addresses: huong.ha@uon.edu.au (Huong Ha), james.welsh@newcastle.edu.au (James S. Welsh), cristian.rojas@ee.kth.se (Cristian R. Rojas), bo.wahlberg@ee.kth.se (Bo Wahlberg).

error.

The paper is organized as follows. Section 2 formulates the problem. Section 3 provides definitions and properties required for the later analysis. The general bound for the SPARSEVA estimation error is then developed in Section 4. In Section 5, we apply the general bound to the special case when the model is cast in a linear regression framework. Section 6 illustrates the developed bound by numerical simulation. Finally, Section 7 provides conclusions.

1.1 Notation

In this paper, we will use the following notation:

- $f(x|0, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp(-x^2/2\sigma^2)$ denotes the probability density function (pdf) of the Normal distribution $\mathcal{N}(0, \sigma^2)$.
- $\chi_\beta^2(N)$ denotes the value that $P(X < \chi_\beta^2(N)) = 1 - \beta$, where X is Chi square distributed with N degrees of freedom.

2 Problem Formulation

Let $Z_1^N = \{Z_1, \dots, Z_N\} \in \mathcal{Z}^N$ denote N identically distributed observations with marginal distribution \mathbb{P} in $\mathcal{Z} \subseteq \mathbb{R}^k$. $\mathcal{L} : \mathbb{R}^n \times \mathcal{Z}^N \rightarrow \mathbb{R}$ denotes a convex and differentiable cost function. Let $\theta^* \in \arg\min_{\theta \in \mathbb{R}^n} \bar{\mathcal{L}}(\theta)$ be a minimizer of the population risk $\bar{\mathcal{L}}(\theta) = \mathbb{E}_{Z_1^N}[\mathcal{L}(\theta; Z_1^N)]$.

The task here is to estimate the unknown parameter θ^* from the data Z_1^N . A well known approach to this problem is to use a regularization technique, i.e., to solve the following convex optimization problem,

$$\hat{\theta}_{\lambda_N} \in \arg \min_{\theta \in \mathbb{R}^n} \{ \mathcal{L}(\theta; Z_1^N) + \lambda_N \mathcal{R}(\theta) \}, \quad (1)$$

where $\lambda_N > 0$ is a user-defined regularization parameter and $\mathcal{R} : \mathbb{R}^n \rightarrow \mathbb{R}^+$ is a norm.

A difficulty when estimating the parameter θ^* using the above regularization technique is that one needs to find the regularization parameter λ_N . The traditional method to choose λ_N is to use cross validation, i.e., to estimate the parameter θ^* with different values of λ_N , then select the value of λ_N that provides the best fit to the validation data. This cross validation method is quite time consuming and very dependent on the data. Here we are specifically interested in the SPARSEVA (SPARSe Estimation based on a Validation criterion) framework, suggested in [16] and [17], which provides automatic tuning of the regularization parameters. Utilizing the SPARSEVA framework, an estimate of θ^* can

be computed using the following convex optimization problem:

$$\begin{aligned} \hat{\theta}_{\varepsilon_N} \in \arg \min_{\theta \in \mathbb{R}^n} \quad & \mathcal{R}(\theta) \\ \text{s.t.} \quad & \mathcal{L}(\theta; Z_1^N) \leq \mathcal{L}(\hat{\theta}_{NR}; Z_1^N)(1 + \varepsilon_N), \end{aligned} \quad (2)$$

where $\varepsilon_N > 0$ is the regularization parameter and $\hat{\theta}_{NR}$ is the “non-regularized” estimate obtained from minimizing the cost function $\mathcal{L}(\theta; Z_1^N)$, i.e.

$$\hat{\theta}_{NR} \in \arg \min_{\theta \in \mathbb{R}^n} \mathcal{L}(\theta; Z_1^N). \quad (3)$$

It can be shown [17] that (1) and (2) are *equivalent* in the sense that there exists a bijection between λ_N and ε_N such that both estimators coincide. However, as discussed in [17, Section V.D], that bijection is data-dependent and it does not seem possible to derive an explicit expression for it. The advantage of the SPARSEVA framework, with respect to (1), is that there are some natural choices of the regularization parameter ε_N based the chosen validation criterion. For example, as suggested in [16] [17], ε_N can be chosen as $2n/N$ (Akaike Information Criterion (AIC)), $n \log(N)/N$ (Bayesian Information Criterion (BIC)); or as suggested in [8], n/N (Prediction Error Criterion).

For the traditional regularization method described in (1), [12] recently developed an upper bound on the estimation error between the estimate $\hat{\theta}_{\lambda_N}$ and the unknown parameter vector θ^* . This bound is a function of some constants related to the nature of the data, the regularization parameter λ_N , the cost function \mathcal{L} and the data length N . The beauty of this bound is that it quantifies the relationship between the estimation error and the finite data length N . Through this relationship, it is easy to confirm most of the properties of the estimate $\hat{\theta}_{\lambda_N}$ in the asymptotic scenario, i.e. $N \rightarrow \infty$, which were developed in the literature some time ago ([9], [11]).

Inspired by [12], our goal is to derive a similar bound for the SPARSEVA estimate $\hat{\theta}_{\varepsilon_N}$, i.e., we want to know how much the SPARSEVA estimate $\hat{\theta}_{\varepsilon_N}$ differs from the true parameter θ^* when the data sample size N is finite. Note that the notation and techniques used in this paper are similar to [12]; however, in [12], the convex optimization problem is posed in the traditional regularization framework (1), while in this paper, the optimization problem is based on the SPARSEVA regularization (2).

3 Definitions and Properties of the Norm $\mathcal{R}(\theta)$ and the Cost Function $\mathcal{L}(\theta)$

In this section, we provide descriptions of some definitions and properties of the norm $\mathcal{R}(\theta)$ and the cost function $\mathcal{L}(\theta; Z_1^N)$, needed to establish an upper bound on the estimation error. Note that we only provide a brief summary such that the research described in this paper can be understood. Readers can find a more detailed discussion in [12].

3.1 Decomposability of a Norm

Let us consider a pair of arbitrary linear subspaces of \mathbb{R}^n , $(\mathcal{M}, \overline{\mathcal{M}})$, such that $\mathcal{M} \subseteq \overline{\mathcal{M}}$. The orthogonal complement of the space $\overline{\mathcal{M}}$ is then defined as,

$$\overline{\mathcal{M}}^\perp = \{v \in \mathbb{R}^n \mid \langle u, v \rangle = 0 \text{ for all } u \in \overline{\mathcal{M}}\},$$

where $\langle \cdot, \cdot \rangle$ is the inner product that maps $\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$.

The norm \mathcal{R} is said to be *decomposable* with respect to $(\mathcal{M}, \overline{\mathcal{M}}^\perp)$ if

$$\mathcal{R}(\theta + \gamma) = \mathcal{R}(\theta) + \mathcal{R}(\gamma) \quad (4)$$

for all $\theta \in \mathcal{M}$ and $\gamma \in \overline{\mathcal{M}}^\perp$.

There are many combinations of norms and vector spaces that satisfy this property (cf. [12]). An example is the l_1 norm and the sparse vector space defined (5). For any subset $S \subseteq \{1, 2, \dots, n\}$ with cardinality s , define the model subspace \mathcal{M} as,

$$\mathcal{M}(S) = \{\theta \in \mathbb{R}^n \mid \theta_j = 0 \text{ for all } j \notin S\}. \quad (5)$$

Now if we define $\overline{\mathcal{M}}(S) = \mathcal{M}(S)$, then the orthogonal complement $\overline{\mathcal{M}}(S)$, with respect to the Euclidean inner product, can be computed as follows,

$$\overline{\mathcal{M}}^\perp(S) = \{\gamma \in \mathbb{R}^n \mid \gamma_j = 0 \text{ for all } j \in S\}.$$

As shown in [12], the l_1 -norm is decomposable with respect to the pair $(\mathcal{M}(S), \overline{\mathcal{M}}^\perp(S))$.

3.2 Dual Norm

For a given inner product $\langle \cdot, \cdot \rangle$, the dual of the norm \mathcal{R} is defined by,

$$\mathcal{R}^*(v) = \sup_{u \in \mathbb{R}^n \setminus \{0\}} \frac{\langle u, v \rangle}{\mathcal{R}(u)} = \sup_{\mathcal{R}(u) \leq 1} \langle u, v \rangle, \quad (6)$$

where \sup is the supremum operator.

Based on the above definition, one can easily see that the dual of the l_1 norm, with respect to the Euclidean inner product, is the l_∞ norm [12].

3.3 Strong Convexity

A twice differentiable function $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$ is *strongly convex* on \mathbb{R}^n when there exists an $m > 0$ such that its Hessian $\nabla^2 \mathcal{L}(\theta)$ satisfies,

$$\nabla^2 \mathcal{L}(\theta) \succeq mI, \quad (7)$$

for all $\theta \in \mathbb{R}^n$ [3]. This is equivalent to the statement that the minimum eigenvalue of $\nabla^2 \mathcal{L}(\theta)$ is not smaller than m for all $\theta \in \mathbb{R}^n$.

An interesting consequence of the strong convexity property in (7) is that for all $\theta, \Delta \in \mathbb{R}^n$, we have,

$$\mathcal{L}(\theta + \Delta) \geq \mathcal{L}(\theta) + \nabla \mathcal{L}(\theta)^T \Delta + \frac{m}{2} \|\Delta\|_2^2. \quad (8)$$

The inequality in (8) has a geometric interpretation in that the graph of the function $\mathcal{L}(\theta)$ has a positive curvature at any $\theta \in \mathbb{R}^n$. The term $m/2$ for the largest m satisfying (7) is typically known as the *curvature* of $\mathcal{L}(\theta)$.

3.4 Subspace Compatibility Constant

For a given norm \mathcal{R} and an error norm $\|\cdot\|$, the *subspace compatibility constant* of a subspace $\mathcal{M} \subseteq \mathbb{R}^n$ with respect to the pair $(\mathcal{R}, \|\cdot\|)$ is defined as,

$$\Psi(\mathcal{M}) = \sup_{u \in \mathcal{M} \setminus \{0\}} \frac{\mathcal{R}(u)}{\|u\|}. \quad (9)$$

This quantity measures how well the norm \mathcal{R} is compatible with the error norm $\|\cdot\|$ over the subspace \mathcal{M} . As shown in [12], when \mathcal{M} is \mathbb{R}^s , the regularized norm \mathcal{R} is the l_1 norm, and the error norm is the l_2 norm, then the subspace compatibility constant is $\Psi(\mathcal{M}) = \sqrt{s}$. Notice also that $\Psi(\mathcal{M})$ is finite, due to the equivalence of finite dimensional norms.

3.5 Projection Operator

The projection of a vector u onto a space \mathcal{M} , with respect to the Euclidean norm, is defined by the following,

$$\Pi_{\mathcal{M}}(u) = \arg \min_{v \in \mathcal{M}} \|u - v\|_2. \quad (10)$$

In the sequel, to simplify the notation, we will write $u_{\mathcal{M}}$ to denote $\Pi_{\mathcal{M}}(u)$.

4 Analysis of the Regularization Technique using the SPARSEVA

In this section, we apply the properties described in Section 3 to derive an upper bound on the error between the SPARSEVA estimate $\hat{\theta}_{\epsilon_N}$ and the unknown parameter θ^* . This upper bound is described in the following theorem.

Theorem 4.1 Assume \mathcal{R} is a norm and is decomposable with respect to the subspace pair $(\mathcal{M}, \overline{\mathcal{M}}^\perp)$ and the cost function $\mathcal{L}(\theta)$ is differentiable and strongly convex with curvature κ_L . Consider the SPARSEVA problem in (2), then the following properties hold:

- i. When $\epsilon_N > 0$, there exists a Lagrange multiplier, $\lambda_N = \lambda_{\epsilon_N}$, such that (1) and (2) have the same solution.

ii. Any optimal solution $\hat{\theta}_{\varepsilon_N} \neq 0$ of the SPARSEVA problem (2) satisfies the following inequalities:

- If ε_N is chosen such that

$$\lambda_{\varepsilon_N} \leq 1/\mathcal{R}^*(\nabla \mathcal{L}(\theta^*)),$$

then

$$\|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2 \leq \frac{4}{\kappa_L^2 \lambda_{\varepsilon_N}^2} \Psi^2(\overline{\mathcal{M}}) + \frac{4}{\kappa_L \lambda_{\varepsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \quad (11)$$

- If ε_N is chosen such that

$$\lambda_{\varepsilon_N} > 1/\mathcal{R}^*(\nabla \mathcal{L}(\theta^*)),$$

then

$$\begin{aligned} \|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2 &\leq \frac{2}{\kappa_L^2} \{\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2 \Psi^2(\overline{\mathcal{M}}) \\ &\quad + \frac{8}{\kappa_L^2} \{\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2 \Psi^2(\overline{\mathcal{M}}^\perp) \\ &\quad + \frac{4}{\kappa_L \lambda_{\varepsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (12)$$

Proof. See the Appendix (Section A.2). \square

Remark 4.1 Note that Theorem 4.1 is intended to provide an upper bound on the estimation error for the general SPARSEVA problem (2). At this stage, it is hard to evaluate, or quantify, the value on the right hand side of the inequalities (11) and (12) as they still contain the term λ_{ε_N} and other abstract terms. However, in the later sections of this paper, from this general upper bound, we will provide bounds on the estimation errors for some specific scenarios.

Remark 4.2 The bound in Theorem 4.1 is actually a family of bounds. For each choice of the pair of subspaces $(\mathcal{M}, \overline{\mathcal{M}}^\perp)$, there is one bound for the estimation error. Hence, in the usual sense, to apply Theorem 4.1 for any specific scenario, the goal is to choose \mathcal{M} and $\overline{\mathcal{M}}^\perp$ to obtain an optimal rate of the bound.

5 An Upper Bound for Sparse Regression

In this section, we illustrate how to apply Theorem 4.1 to derive an upper bound of the error between the SPARSEVA estimate $\hat{\theta}_{\varepsilon_N}$ and the true parameter θ^* for the following linear regression model,

$$Y_N = \Phi_N^T \theta^* + e, \quad (13)$$

where $\theta^* \in \mathbb{R}^n$ is the unknown parameter that is required to be estimated; $e \in \mathbb{R}^N$ is the disturbance noise; $\Phi_N \in \mathbb{R}^{n \times N}$ is

the regression matrix and $Y_N \in \mathbb{R}^N$ is the output vector. Here, we make the following assumption on the true parameter θ^* ,

Assumption 5.1 The true parameter θ^* is “weakly” sparse, i.e. $\theta^* \in \mathbb{B}_q(R_q)$, where,

$$\mathbb{B}_q(R_q) := \left\{ \theta \in \mathbb{R}^n \left| \sum_{i=1}^p |\theta_i|^q \leq R_q \right. \right\}, \quad (14)$$

with $q \in [0, 1]$ being a constant.

Using the SPARSEVA framework in (2) with \mathcal{R} chosen as the l_1 norm and the cost function $\mathcal{L}(\theta)$ chosen as,

$$\mathcal{L}(\theta) = \frac{1}{2N} \|Y_N - \Phi_N^T \theta\|_2^2, \quad (15)$$

then an estimate of θ^* in (13) can be found by solving the following problem,

$$\begin{aligned} \hat{\theta}_{\varepsilon_N} &\in \arg \min_{\theta \in \mathbb{R}^n} \|\theta\|_1 \\ \text{s.t.} \quad &\mathcal{L}(\theta) \leq \mathcal{L}(\hat{\theta}_{NR})(1 + \varepsilon_N), \end{aligned} \quad (16)$$

with $\hat{\theta}_{NR} = (\Phi_N \Phi_N^T)^{-1} \Phi_N Y_N$ and $\varepsilon_N > 0$ being the user-defined regularization parameter. Now ε_N can be chosen as either $2n/N$ or $\log(N)n/N$ as suggested in [16]; or n/N as suggested in [8].

Remark 5.1 Note that the sparse regression problem is very common in system identification and is often used to obtain a low order linear model by regularization.

Remark 5.2 For Assumption 5.1, note that when $q = 0$, under the convention that $0^0 = 0$, the set in (14) corresponds to an exact sparsity set, where all the elements belonging to the set have at most R_0 non-zero entries. Generally, for $q \in (0, 1]$, the set $\mathbb{B}_q(R_q)$ forces the ordered absolute values of θ^* to decay with a certain rate.

5.1 An Analysis on the Strong Convexity Property and the Curvature of the l_2 norm Cost Function

Consider the convex optimization problem in (16), the Hessian matrix of the cost function $\mathcal{L}(\theta)$ is computed as,

$$\nabla^2 \mathcal{L}(\theta) = \frac{1}{N} \Phi_N \Phi_N^T.$$

To prove that $\mathcal{L}(\theta)$ is strongly convex, we need to prove,

$$\exists \kappa_L > 0 \quad \text{s.t.} \quad \frac{1}{N} \Phi_N \Phi_N^T \succeq 2\kappa_L I. \quad (17)$$

We see that the requirement in (17) coincides with the requirement of persistent excitation of the input signal in a

system identification problem. If an experiment is well-designed, then the input signal $u(t)$ needs to be persistently exciting of order n , i.e., the matrix $\Phi_N \Phi_N^T$ is a positive definite matrix. This means that the condition in (17) is always satisfied for any linear regression problem derived from a well posed system identification problem. This means that for any choice of the regression matrix Φ_N that satisfies the persistent excitation condition, there exists a positive curvature κ_L of the cost function $\mathcal{L}(\theta)$.

Consider $\Phi_N \in \mathbb{R}^{n \times N}$ to be a matrix where each row $\Phi_{N,j}$ is sampled from a Normal distribution of zero mean and covariance matrix $\Sigma \in \mathbb{R}^{N \times N}$, i.e., $\Phi_{N,j} \sim \mathcal{N}(0, \Sigma)$, $\forall j = 1, \dots, n$. We then denote the distribution of the smallest eigenvalue of $N^{-1} \Phi_N \Phi_N^T$ to be $P(x|\Sigma, N, n)$, means that given a probability $1 - \alpha$, $0 \leq \alpha \leq 1$, there exists a value w_{\min} such that $N^{-1} \Phi_N \Phi_N^T \succeq w_{\min} I$, for any matrix Φ_N constructed following the above assumption. Then the global curvature κ , i.e. the curvature that satisfies (17) for any regression matrix Φ_N , can be expressed as $(1/2)w_{\min}$. For the rest of the paper, we will denote by κ_α lower bound on the global curvature κ with probability $1 - \alpha$, $0 \leq \alpha \leq 1$.

5.2 Assumptions

For the linear regression in (13), the following assumptions are made:

Assumption 5.2 *The rows $\Phi_{N,j}, j = 1, \dots, n$ of the regressor matrix Φ_N are distributed as $\Phi_{N,j} \sim \mathcal{N}(0, \Sigma)$, where $\Sigma \in \mathbb{R}^{N \times N}$ is a constant, symmetric, positive definite matrix.*

Note that an obvious practical case where Assumption 5.2 is satisfied is when the model is FIR and the input signal being white noise or coloured noise.

Assumption 5.3 *The noise vector $e \in \mathbb{R}^N$ is Gaussian with i.i.d. $\mathcal{N}(0, \sigma_e^2)$ entries.*¹

5.3 Developing the Upper Bound

The following theorem provides an upper bound on the estimation error $\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2$ for the optimization problem in (16) in the case of weakly sparse estimates.

Theorem 5.1 *Suppose Assumptions 5.2, 5.3 and 5.1 hold, when N is large, then with probability $(1 - \alpha)(1 - 4n\beta)$ ($0 \leq \alpha \leq 1$, $0 \leq \beta \leq 1$), if $\hat{\theta}_{\epsilon_N} \neq 0$ we have the following inequality*

$$\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2^2 \leq \max(a_1, a_2), \quad (18)$$

¹ The assumption of Gaussian noise is fairly standard in system identification. However, this assumption can be relaxed to ‘sub-Gaussian’ noise (i.e., when the tails of the noise distribution decay like $e^{-\alpha x^2}$) at the expense of longer derivations.

where

$$\begin{aligned} a_1 &= \frac{8n_\eta \sigma_e^2 s_{\max} \chi_\beta^2 (N - n)(1 + \epsilon_N) \ln(2/\beta)}{\kappa_\alpha^2 N^2} \\ &\quad + \frac{\sqrt{32 \sigma_e^2 s_{\max} \chi_\beta^2 (N - n)(1 + \epsilon_N) \ln(2/\beta)}}{\kappa_\alpha N} \|\theta_{[n_\eta+1:n]}^*\|_1, \\ a_2 &= \frac{(16n - 12n_\eta) \sigma_e^2 \chi_\beta^2 (\Sigma, I) \ln(2/\beta)}{\kappa_\alpha^2 N^2} \\ &\quad + \frac{\sqrt{32 \sigma_e^2 s_{\max} \chi_\beta^2 (N - n)(1 + \epsilon_N) \ln(2/\beta)}}{\kappa_\alpha N} \|\theta_{[n_\eta+1:n]}^*\|_1. \end{aligned}$$

where κ_α is a lower bound on the curvature of the regression matrix (i.e., half the smallest eigenvalue of $N^{-1} \Phi_N^T \Phi_N$) with probability $1 - \alpha$, n_η is any integer between 1 and n , $\theta_{[n_\eta+1:n]}^*$ is the vector formed from the $n - n_\eta$ smallest (in magnitude) entries of θ^* , and s_{\max} is the maximum singular value of the matrix Σ .

Proof. This proof relies on three preliminary results introduced in Appendix A.3. For an integer $n_\eta \in \{1, \dots, n\}$, define S_η as the set of the indices of the n_η largest (in magnitude) entries of θ^* , and its complementary set S_η^c as

$$S_\eta^c = \{1, 2, \dots, n\} \setminus S_\eta; \quad (19)$$

with the corresponding subspaces $\mathcal{M}(S_\eta)$ and $\mathcal{M}^\perp(S_\eta)$ as,

$$\begin{aligned} \mathcal{M}(S_\eta) &= \{\theta \in \mathbb{R}^n \mid \theta_j = 0 \ \forall j \notin S_\eta\}, \\ \mathcal{M}^\perp(S_\eta) &= \{\gamma \in \mathbb{R}^n \mid \gamma_j = 0 \ \forall j \in S_\eta\}. \end{aligned} \quad (20)$$

Using the definition of the subspace compatibility constant described in Section 3, we have,

$$\begin{aligned} \Psi^2(\mathcal{M}(S_\eta)) &= |S_\eta| = n_\eta, \\ \Psi^2(\mathcal{M}^\perp(S_\eta)) &= |S_\eta^c| = n - n_\eta. \end{aligned} \quad (21)$$

where $|S|$ denotes the cardinality of S .

Now, for Theorem 4.1 to generate an upper bound for the problem (16), we need to establish an upper bound on $\|\theta_{\mathcal{M}^\perp(S_\eta)}^*\|_1$. Based on the definition of the subspace $\mathcal{M}^\perp(S_\eta)$, we have,

$$\|\theta_{\mathcal{M}^\perp(S_\eta)}^*\|_1 = \|\theta_{[n_\eta+1:n]}^*\|_1, \quad (22)$$

where $\theta_{[n_\eta+1:n]}^*$ denotes the vector formed from the $n - n_\eta$ smallest (in magnitude) entries of θ^* . Define κ_α as a lower bound on the global curvature of the regression matrix, i.e. half the smallest eigenvalue of $\Phi_N^T \Phi_N$, with probability $1 - \alpha$, $0 \leq \alpha \leq 1$. Substituting the results of Propositions A.1-A.3 from Appendix A.3, (21) and (22) into the bound in

Theorem 4.1, then with n_η being any integer between 1 and n , we have the following bounds:

- If ε_N is chosen such that

$$\lambda_{\varepsilon_N} \leq 1/\mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) = 1/\|\nabla \mathcal{L}(\theta^*)\|_\infty,$$

then, with probability at least $(1 - \alpha)(1 - 2n\beta)$,

$$\begin{aligned} & \|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2 \\ & \leq \frac{8n_\eta \sigma_e^2 s_{\max} \chi_\beta^2(N - n)(1 + \varepsilon_N) \ln(2/\beta)}{\kappa_\alpha^2 N^2} \\ & \quad + \frac{\sqrt{32\sigma_e^2 s_{\max} \chi_\beta^2(N - n)(1 + \varepsilon_N) \ln(2/\beta)}}{\kappa_\alpha N} \|\theta_{[n_\eta+1:n]}^*\|_1. \end{aligned}$$

- If ε_N is chosen such that

$$\lambda_{\varepsilon_N} > 1/\mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) = 1/\|\nabla \mathcal{L}(\theta^*)\|_\infty,$$

then, with probability at least $(1 - \alpha)(1 - 4n\beta)$,

$$\begin{aligned} & \|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2 \\ & \leq \frac{(16n - 12n_\eta) \sigma_e^2 \chi_\beta^2(\Sigma, I) \ln(2/\beta)}{\kappa_\alpha^2 N^2} \\ & \quad + \frac{\sqrt{32\sigma_e^2 s_{\max} \chi_\beta^2(N - n)(1 + \varepsilon_N) \ln(2/\beta)}}{\kappa_\alpha N} \|\theta_{[n_\eta+1:n]}^*\|_1. \end{aligned}$$

Therefore, for n_η being any integer between 1 and n , with probability at least $(1 - \alpha)(1 - 4n\beta)$, we have

$$\|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2 \leq \max(a_1, a_2), \quad (23)$$

where

$$\begin{aligned} a_1 &= \frac{8n_\eta \sigma_e^2 s_{\max} \chi_\beta^2(N - n)(1 + \varepsilon_N) \ln(2/\beta)}{\kappa_\alpha^2 N^2} \\ & \quad + \frac{\sqrt{32\sigma_e^2 s_{\max} \chi_\beta^2(N - n)(1 + \varepsilon_N) \ln(2/\beta)}}{\kappa_\alpha N} \|\theta_{[n_\eta+1:n]}^*\|_1, \\ a_2 &= \frac{(16n - 12n_\eta) \sigma_e^2 \chi_\beta^2(\Sigma, I) \ln(2/\beta)}{\kappa_\alpha^2 N^2} \\ & \quad + \frac{\sqrt{32\sigma_e^2 s_{\max} \chi_\beta^2(N - n)(1 + \varepsilon_N) \ln(2/\beta)}}{\kappa_\alpha N} \|\theta_{[n_\eta+1:n]}^*\|_1. \end{aligned}$$

□

Remark 5.3 The bound in Theorem 5.1 is also a family of bounds, one for each value of n_η .

Remark 5.4 When $\Sigma = \sigma_u I$, i.e. the model is FIR and the input $u(t)$ is white noise, then $s_{\max} = \sigma_u$ and the generalized Chi square distribution $\chi^2(\Sigma, I)$ becomes the Chi square distribution $\sigma_u \chi^2(N)$.

Remark 5.5 Note that the developed bound in Theorem 5.1 depends on the true parameter θ^* , which is unknown but constant. Using a similar proof as in Proposition 2.3 of [7], we can derive under Assumption 5.1 an upper bound for the term $\|\theta_{[n_\eta+1:n]}^*\|_1$. Specifically, we have,

$$\|\theta_{[n_\eta+1:n]}^*\|_1 = \sum_{i=n_\eta+1}^n |\theta_{[i]}^*| = \sum_{i=n_\eta+1}^n |\theta_{[i]}^*|^{1-q} |\theta_{[i]}^*|^q$$

Since S_η is the set of the indices of the n_η largest (in magnitude) entries of θ^* , i.e. $|\theta_{[i]}^*| \leq |\theta_{n_\eta}^*|$, $\forall i = n_\eta + 1, \dots, n$, hence,

$$\|\theta_{[n_\eta+1:n]}^*\|_1 \leq |\theta_{n_\eta}^*|^{1-q} \sum_{i=n_\eta+1}^n |\theta_{[i]}^*|^q$$

Using the same argument, we have,

$$|\theta_{n_\eta}^*|^{1-q} = \left(\frac{1}{n_\eta} \sum_{i=1}^{n_\eta} |\theta_{[i]}^*|^q \right)^{(1-q)/q} \leq \left(\frac{1}{n_\eta} \sum_{i=1}^{n_\eta} |\theta_{[i]}^*|^q \right)^{(1-q)/q}.$$

Therefore,

$$\begin{aligned} \|\theta_{[n_\eta+1:n]}^*\|_1 & \leq \left(\frac{1}{n_\eta} \sum_{i=1}^{n_\eta} |\theta_{[i]}^*|^q \right)^{(1-q)/q} \sum_{i=n_\eta+1}^n |\theta_{[i]}^*|^q \\ & \leq \left(\frac{1}{n_\eta} \sum_{i=1}^{n_\eta} |\theta_{[i]}^*|^q \right)^{(1-q)/q} \sum_{i=1}^n |\theta_{[i]}^*|^q \\ & \leq \left(\frac{1}{n_\eta} \right)^{1/q-1} \|\theta^*\|_q^{1-q} \|\theta^*\|_q^q \\ & \leq (n_\eta)^{1-1/q} \|\theta^*\|_q \\ & \leq (n_\eta)^{1-1/q} (R_q)^{1/q}. \end{aligned}$$

This means we can always place an upper bound on the term $\|\theta_{[n_\eta+1:n]}^*\|_1$ by a known constant which depends on the nature of the true parameter θ^* . Therefore, from Theorem 5.1, we can see that the estimation error $\|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2 = O_p(N^{-1/2})$ [17]. This confirms the result in [17], that in the asymptotic case, when $\varepsilon_N > 0$, the SPARSEVA estimate $\hat{\theta}_{\varepsilon_N}$ converges to the true parameter θ^* .

6 Numerical Evaluation

In this section, numerical examples are presented to illustrate the bound $\|\hat{\theta}_{\varepsilon_N} - \theta^*\|_2^2$ as stated in Theorem 5.1. In Section 6.1, we consider the case when the input is Gaussian white noise whilst in Section 6.2, the input is a correlated signal with zero mean.

6.1 Gaussian White Noise Input

In this section, a random discrete time system with a random model order between 1 and 10 is generated using the command *drss* from Matlab. The system has poles with magnitude less than 0.9. Gaussian white noise is added to the system output to give different levels of SNR, e.g. 30dB, 20dB and 10dB. For each noise level, 50 different input excitation signals (Gaussian white noise with variance 1) and output noise realizations are generated. For each set of input and output data, the system parameters are estimated using a different sample size, i.e., $N = [450, 1000, 5000, 10000, 50000, 100000]$.

The FIR model structure is used here in order to construct the SPARSEVA problem (16). The number of parameters n of the FIR model is set to be 35. The regularization parameter ϵ_N is chosen as n/N [8].

We then compute the upper bound of $\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2$ using (18) with different values of n_η , i.e. $n_\eta = [10, 15, 25]$. The probability parameters α and β are chosen to be 0.02 and 0.001 respectively. Related to the computation of the universal constant κ corresponding to the distribution $\mathcal{N}(0, \Sigma)$, note that, in reality, it is very difficult to compute its exact distribution $P(x|\Sigma, N, n)$, hence, here we use an empirical method to compute the distribution $P(x|\Sigma, N, n)$. The idea is to generate a large number of random matrices Φ_N , compute the smallest eigenvalue of $N^{-1}\Phi_N\Phi_N^T$, and then build a histogram of these values, which is an approximation of $P(x|\Sigma, N, n)$. Then we compute the value of w_{\min} to ensure the inequality $N^{-1}\Phi_N\Phi_N^T \succeq w_{\min}I$ occurs with probability $1 - \alpha$. Finally, κ_α is computed using the formula $\kappa_\alpha = (1/2)w_{\min}$.

With the setting described above, the probability of the upper bound being correct is $(1 - \alpha)(1 - 4n\beta) = 0.84$. This upper bound will be compared with $\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2$. Note that we plot both the upper bound and the true estimation errors on a logarithmic scale.

Plots of the estimation error versus the data length N with different noise levels are displayed in Figures 1 to 3. In Figures 1 to 3, the red lines are the true estimation errors from 50 estimates using the SPARSEVA framework. The magenta, blue and cyan lines are the upper bounds developed in Theorem 5.1, which correspond to $n_\eta = [10, 15, 25]$, respectively. We can see that the plots confirm the bound developed in Theorem 5.1 for all noise levels. When N becomes large, the estimation error and the corresponding upper bound become smaller. When N goes to infinity, the estimation error will tend to 0. Note that the bounds are slightly different for the chosen values of n_η , however, not significantly. As can be seen, the bounds are relatively insensitive to the choice of n_η .

In addition, we plot another graph, shown in Fig. 4, to compare the proposed upper bound and the true estimation errors corresponding to different value of ϵ_N , i.e. $\frac{n}{N}$ (PEC),

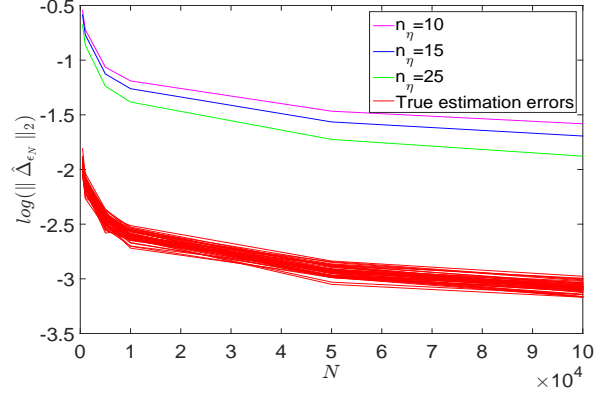


Fig. 1. Plot of the estimation error for a SNR=30dB and a Gaussian white input signal.

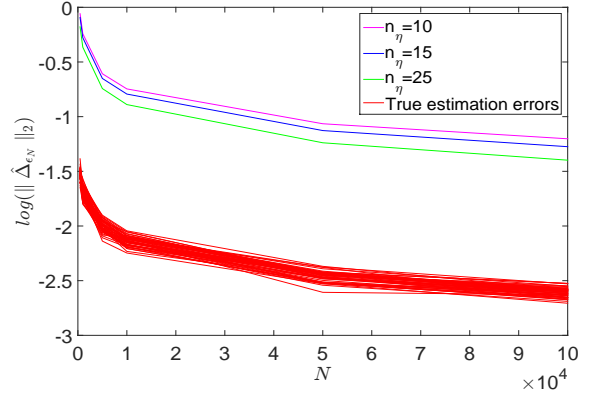


Fig. 2. Plot of the estimation error for a SNR=20dB and a Gaussian white input signal.

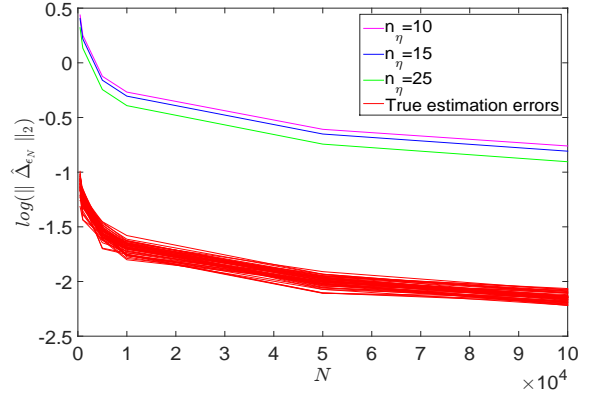


Fig. 3. Plot of the estimation error for a SNR=10dB and a Gaussian white input signal.

$\frac{2n}{N}$ (AIC) and $\frac{\log(N)n}{N}$ (BIC). The blue lines are the upper bounds developed in Theorem 5.1, which correspond to $n_\eta = 25$, with the three different values of ϵ_N . The magenta (BIC), green (AIC) and red (PEC) lines are the true estimation errors from 50 estimates (for each value of ϵ_N) using

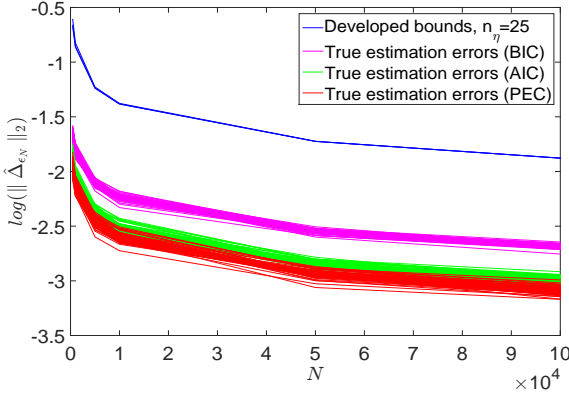


Fig. 4. Plot of the proposed bound and the true estimation errors corresponding to different choices of ϵ_N , for a SNR=30dB and a white Gaussian input signal (magenta BIC, green AIC and red PEC).

the SPARSEVA framework. We can see that the plot again confirms the validity of the proposed upper bound for all choices of ϵ_N . Note that the upper bound is not extremely tight, it is quite conservative, however, it is the price to usually pay for finite sample bounds with a general SPARSEVA setting, i.e. the regularized parameter ϵ_N can be any positive value. When ϵ_N is larger, the upper bound will be closer to the true estimate error.

6.2 Coloured Noise Input

In this section, a random discrete time system with a random model order between 1 and 10 is generated using the command *drss* from Matlab. The system has poles with magnitude less than 0.9. White noise is added to the system output with different levels of SNR, e.g. 30dB, 20dB and 10dB. For each noise level, 50 different input excitation signals and output noise realizations are generated. For each set of input and output data, the system parameters are estimated using different sample sizes, i.e. $N = [450, 1000, 5000, 10000, 50000]$.

Here, the input signal is generated by filtering a zero mean Gaussian white noise with unit variance through the filter,

$$F_u(q) = \frac{0.9798}{1 - 0.2q^{-1}}.$$

Due to this filtering, the covariance matrix of the regression matrix distribution will not be of a diagonal form. Note that this is a completely different scenario to that in Section 6.1.

The FIR model structure is used here in order to construct the linear regression for the SPARSEVA problem (16). The number of parameters n of the FIR model is set to be 35. The regularization parameter, ϵ_N , is chosen as n/N [8].

We then compute the upper bound of $\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2$ using (18) with different values of n_η , i.e. $n_\eta = [10, 15, 25]$. The

probability parameters α and β are chosen to be 0.02 and 0.001 respectively. With this setting, the probability of the upper bound being correct is $(1 - \alpha)(1 - 4n\beta) = 0.84$. This upper bound will be compared with $\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2$.

Plots of the upper bound as stated in Theorem 5.1 and the true estimation error $\|\hat{\theta}_{\epsilon_N} - \theta^*\|_2$ are displayed in Figures 5 to 7. In Figures 5 to 7, the red lines are the true estimation errors from 50 estimates using the SPARSEVA framework. The magenta, blue and cyan lines are the upper bounds developed in Theorem 5.1, which correspond to $n_\eta = [10, 15, 25]$ respectively. We can see that the plots confirmed the bound developed in Theorem 5.1 for all noise levels. When N becomes large, the estimation error and the corresponding upper bound become smaller. When N goes to infinity, the estimation error will tend to 0. Note that the bounds are slightly different for the chosen values of n_η , however, not significantly. As can be seen, the bounds are relatively insensitive to the choice of n_η .

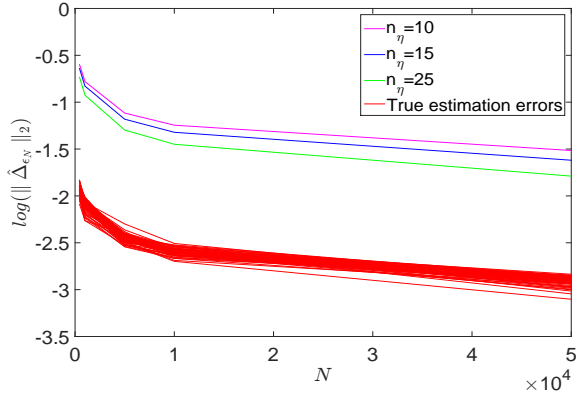


Fig. 5. Plot of the estimation error for SNR=30dB for a coloured input signal.

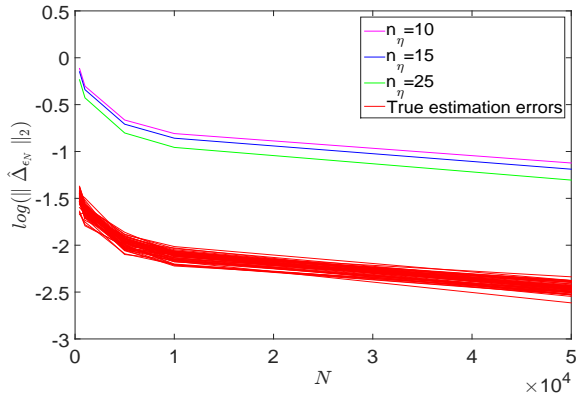


Fig. 6. Plot of the estimation error for SNR=20dB for a coloured input signal.

7 Conclusion

The paper provides an upper bound on the SPARSEVA estimation error in the general case, for any choice of strongly

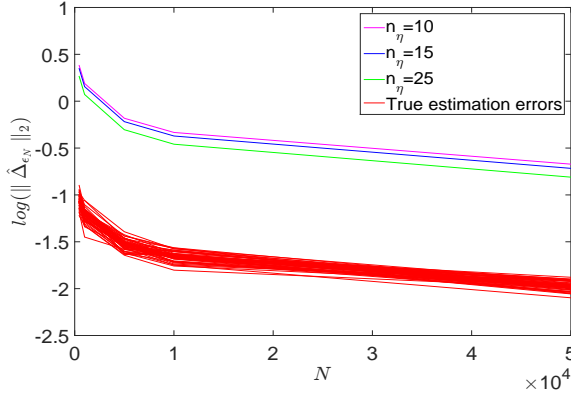


Fig. 7. Plot of the estimation error for SNR=10dB for a coloured input signal.

convex cost function and decomposable norm. We also evaluate the bound for a specific scenario, i.e., a sparse regression estimate problem. Numerical results confirm the validity of the developed bound for different input signals with different output noise levels for different choices of the regularization parameters.

References

- [1] P.J. Bickel, Y. Ritov, and A.B. Tsybakov. Simultaneous analysis of lasso and Dantzig selector. *Annals of Statistics*, 37:1705–1732, 2009.
- [2] J.M. Borwein and Q.J. Zhu. A variational approach to lagrange multipliers. *Journal of Optimization Theory and Applications*, 171(3):727–756, 2016.
- [3] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [4] E. Candes and T. Tao. The Dantzig selector: Statistical estimation when p is much larger than n . *Annals of Statistics*, 35:2313–2351, 2007.
- [5] M. Fazel, H. Hindi, and S. Boyd. A rank minimization heuristic with application to minimum order system approximation. *Proceedings 2001 American Control Conference*, 2001.
- [6] M. Fazel, H. Hindi, and S. Boyd. Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices. *Proceedings 2003 American Control Conference*, 2003.
- [7] S. Foucart and H. Rauhut. *A Mathematical Introduction to Compressive Sensing*. Birkhäuser Basel, 2013.
- [8] H. Ha, J.S. Welsh, N. Blomberg, C.R. Rojas, and B. Wahlberg. Reweighted nuclear norm regularization: A SPARSEVA approach. *Proceedings of the 17th IFAC Symposium on System Identification*, 48(28):1172–1177, 2015.
- [9] J. Huang, J.L. Horowitz, and S. Ma. Asymptotic properties of bridge estimators in sparse high-dimensional regression models. *Annals of Statistics*, 36(2):587–613, 2008.
- [10] A.T. James. Distributions of matrix variates and latent roots derived from normal samples. *The Annals of Mathematical Statistics*, 35(2):475–501, 1964.
- [11] K. Knight and W. Fu. Asymptotics for Lasso-Type Estimators. *Annals of Statistics*, 28(5):1356–1378, 2000.
- [12] S.N. Negahban, P. Ravikumar, M.J. Wainwright, and B. Yu. A Unified Framework for High-Dimensional Analysis of M-Estimators with Decomposable Regularizers. *Statistical Science*, 27(4):538–557, 2012.
- [13] H. Ohlsson. *Regularization for Sparseness and Smoothness – Applications in System Identification and Signal Processing*. PhD thesis, Department of Electrical Engineering, Linköping University, Sweden, 2010.
- [14] M. Osborne, B. Presnell, and B. Turlach. On the LASSO and its dual. *Journal of Computational and Graphical Statistics*, 9(2):319–337, 2000.
- [15] G. Pillonetto, F. Dinuzzo, T. Chen, G. De Nicolao, and L. Ljung. Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50(3):657–682, 2014.
- [16] C.R. Rojas and H. Hjalmarsson. Sparse estimation based on a validation criterion. *2011 50th IEEE Conference on Decision and Control and European Control Conference (CDC-ECC)*, pages 2825–2830, 2011.
- [17] C.R. Rojas, R. Tóth, and H. Hjalmarsson. Sparse estimation of polynomial and rational dynamical models. *IEEE Transactions on Automatic Control*, 59:2962–2977, 2014.
- [18] T. Söderström and P. Stoica. *System Identification*. Prentice Hall, 1989.
- [19] R. Tibshirani. Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- [20] A.N. Tikhonov and V.Y. Arsenin. *Solutions of Ill-Posed Problems*. Winston & Sons, 1977.
- [21] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. In Y.C. Eldar and G. Kutyniok, editors, *Compressed Sensing: Theory and Applications*. Cambridge University Press, 2012.
- [22] M. Zorzi and A. Chiuso. Sparse plus low rank network identification: A nonparametric approach. *Automatica*, 76:355–366, 2017.

A Appendix

A.1 Background knowledge

First, we cite a lemma directly from [12], to enable the proof of Theorem 4.1 to be constructed.

Lemma A.1 *For any norm \mathcal{R} that is decomposable with respect to $(\mathcal{M}, \overline{\mathcal{M}}^\perp)$; and any vectors θ, Δ , we have*

$$\mathcal{R}(\theta + \Delta) - \mathcal{R}(\theta) \geq \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp}) - \mathcal{R}(\Delta_{\mathcal{M}}) - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}), \quad (\text{A.1})$$

Recall that $\Delta_{\overline{\mathcal{M}}^\perp}$ is the Euclidean projection of Δ onto

$\overline{\mathcal{M}}^\perp$ (see Section 3.5), and similarly for the other terms in (A.1).

Proof. See the supplementary material of [12]. \square

We now quote the following lemma from [2] with modification to fit with the notation used in the SPARSEVA problem (2). This lemma helps us to find some important properties related to the SPARSEVA estimate $\hat{\theta}_{\lambda_N}$. Based on these properties, in the next section we can derive the upper bound on the estimation error. Note that for notational simplicity, we will denote $\mathcal{L}(\theta; Z_1^N)$ as $\mathcal{L}(\theta)$.

Lemma A.2 *Consider the convex optimization problem in (2). Then the pair $(\hat{\theta}_{\epsilon_N}, \lambda_{\epsilon_N})$, with $\hat{\theta}_{\epsilon_N} \neq 0$, has the property that $\hat{\theta}_{\epsilon_N}$ is the solution of the problem (2) and λ_{ϵ_N} is the Lagrange multiplier if and only if all of the following hold:*

- (1) $\lambda_{\epsilon_N} \in \mathbb{R}^+$;
- (2) the function $\mathcal{R}(\theta) + \lambda\{\mathcal{L}(\theta) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)\}$ attains its minimum over \mathbb{R}^n at $\hat{\theta}_{\epsilon_N}$; and
- (3) $\mathcal{L}(\hat{\theta}_{\epsilon_N}) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N) = 0$.

Proof. See the proof of Theorem 2.2 in [2]. The third condition in the cited theorem is a complementary slackness condition, which reduces to condition (3) here if $\hat{\theta}_{\epsilon_N} \neq 0$ (cf. [17, Lemma II.2]). \square

A.2 Proof of Theorem 4.1

First we need to prove that there exists a Lagrange multiplier for the SPARSEVA problem (2). We can assume without loss of generality that $\mathcal{L}(\hat{\theta}_{NR}) \neq 0$, since otherwise we can take $\lambda_{\epsilon_N} = 0$. According to [2], the Lagrange multiplier for a convex optimization problem with constraint exists when the Slater condition is satisfied. Specifically, for the SPARSEVA problem (2), the Lagrange multiplier λ_{ϵ_N} exists when there exists a θ_1 such that $\mathcal{L}(\theta_1) < \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)$. If $\epsilon_N > 0$, and $\mathcal{L}(\hat{\theta}_{NR}) \neq 0$, there always exists a parameter vector θ_1 such that $\mathcal{L}(\theta_1) < \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)$ (just take $\theta_1 = \hat{\theta}_{NR}$). Therefore, there exists a Lagrange multiplier λ_{ϵ_N} for the SPARSEVA problem.

Now that we have confirmed the existence of the Lagrange multiplier λ_{ϵ_N} , consider the function $\mathcal{F}(\Delta)$ defined as follows,

$$\mathcal{F}(\Delta) = \mathcal{R}(\theta^* + \Delta) - \mathcal{R}(\theta^*) + \lambda_{\epsilon_N}\{\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*)\}. \quad (\text{A.2})$$

Using the strong convexity condition of $\mathcal{L}(\theta)$,

$$\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) \geq \langle \nabla \mathcal{L}(\theta^*), \Delta \rangle + \kappa_{\mathcal{L}} \|\Delta\|_2^2. \quad (\text{A.3})$$

From (6), we have that

$$|\langle \nabla \mathcal{L}(\theta^*), \Delta \rangle| \leq \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \mathcal{R}(\Delta). \quad (\text{A.4})$$

Next, combining the inequality (A.4) and the triangle inequality, i.e., $\mathcal{R}(\Delta) \leq \mathcal{R}(\Delta_{\overline{\mathcal{M}}}) + \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp})$, we have,

$$\begin{aligned} |\langle \nabla \mathcal{L}(\theta^*), \Delta \rangle| &\leq \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \mathcal{R}(\Delta) \\ &\leq \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) (\mathcal{R}(\Delta_{\overline{\mathcal{M}}}) + \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp})), \end{aligned}$$

therefore,

$$\langle \nabla \mathcal{L}(\theta^*), \Delta \rangle \geq -\mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) (\mathcal{R}(\Delta_{\overline{\mathcal{M}}}) + \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp})). \quad (\text{A.5})$$

Now, combining (A.2), (A.3), (A.5), and Lemma A.1,

$$\begin{aligned} \mathcal{F}(\Delta) &= \mathcal{R}(\theta^* + \Delta) - \mathcal{R}(\theta^*) + \lambda_{\epsilon_N}\{\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*)\} \\ &\geq \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp}) - \mathcal{R}(\Delta_{\overline{\mathcal{M}}}) - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}^*) \\ &\quad + \lambda_{\epsilon_N}\left\{-\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))(\mathcal{R}(\Delta_{\overline{\mathcal{M}}}) + \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp})) + \kappa_{\mathcal{L}} \|\Delta\|_2^2\right\} \\ &\geq \{1 - \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \mathcal{R}(\Delta_{\overline{\mathcal{M}}^\perp}) \\ &\quad - \{1 + \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \mathcal{R}(\Delta_{\overline{\mathcal{M}}}) + \kappa_{\mathcal{L}} \lambda_{\epsilon_N} \|\Delta\|_2^2 \\ &\quad - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (\text{A.6})$$

Notice that, when $\hat{\theta}_{\epsilon_N}$ is the estimate of the SPARSEVA problem, then property 2 in Lemma A.2 states that the function $\mathcal{R}(\theta) + \lambda\{\mathcal{L}(\theta) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)\}$ attains its minimum over \mathbb{R}^n at $\hat{\theta}_{\epsilon_N}$, which means,

$$\begin{aligned} \forall \theta \in \mathcal{R}^n, \mathcal{R}(\hat{\theta}_{\epsilon_N}) + \lambda_{\epsilon_N}(\mathcal{L}(\hat{\theta}_{\epsilon_N}) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)) \\ \leq \mathcal{R}(\theta) + \lambda_{\epsilon_N}(\mathcal{L}(\theta) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)). \end{aligned}$$

Hence,

$$\forall \theta \in \mathcal{R}^n, \mathcal{R}(\hat{\theta}_{\epsilon_N}) - \mathcal{R}(\theta) + \lambda_{\epsilon_N}\{\mathcal{L}(\hat{\theta}_{\epsilon_N}) - \mathcal{L}(\theta)\} \leq 0,$$

or, taking $\theta = \theta^*$ and defining $\hat{\Delta}_{\epsilon_N} := \hat{\theta}_{\epsilon_N} - \theta^*$,

$$\mathcal{F}(\hat{\Delta}_{\epsilon_N}) \leq 0. \quad (\text{A.7})$$

Combining (A.6) with (A.7), we then have,

$$\begin{aligned} 0 &\geq \kappa_{\mathcal{L}} \lambda_{\epsilon_N} \|\hat{\Delta}_{\epsilon_N}\|_2^2 + \{1 - \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \mathcal{R}(\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}^\perp}) \\ &\quad - \{1 + \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \mathcal{R}(\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}}) - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (\text{A.8})$$

Now we consider two cases.

Case 1: $\lambda_{\epsilon_N} \leq 1/\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))$

From (A.8), we have,

$$\begin{aligned} 0 &\geq \kappa_{\mathcal{L}} \lambda_{\epsilon_N} \|\hat{\Delta}_{\epsilon_N}\|_2^2 - \{1 + \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \mathcal{R}(\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}}) \\ &\quad - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (\text{A.9})$$

By the definition of subspace compatibility,

$$\mathcal{R}(\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}}) \leq \Psi(\overline{\mathcal{M}}) \|\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}}\|_2.$$

Now we also have that

$$\|\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}}\|_2 = \|\Pi_{\overline{\mathcal{M}}}(\hat{\Delta}_{\epsilon_N})\|_2 \leq \|\hat{\Delta}_{\epsilon_N}\|_2.$$

Therefore,

$$\mathcal{R}(\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}}) \leq \Psi(\overline{\mathcal{M}}) \|\hat{\Delta}_{\epsilon_N}\|_2. \quad (\text{A.10})$$

Substituting this into (A.9) gives,

$$0 \geq \kappa_L \lambda_{\epsilon_N} \|\hat{\Delta}_{\epsilon_N}\|_2^2 - \{1 + \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \Psi(\overline{\mathcal{M}}) \|\hat{\Delta}_{\epsilon_N}\|_2 - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \quad (\text{A.11})$$

Note that for a quadratic polynomial $f(x) = ax^2 + bx + c$, with $a > 0$, if there exists $x \in \mathbb{R}^+$ that makes $f(x) \leq 0$, then such x must satisfy

$$x \leq \frac{-b + \sqrt{b^2 - 4ac}}{2a}.$$

Since $(A + B)^2 \leq 2A^2 + 2B^2$ for all $A, B \in \mathbb{R}$,

$$x^2 \leq 2 \left[\frac{b^2}{4a^2} + \frac{b^2 - 4ac}{4a^2} \right] = \frac{b^2 - 2ac}{a^2}. \quad (\text{A.12})$$

Applying this inequality to (A.11), we have,

$$\begin{aligned} \|\hat{\Delta}_{\epsilon_N}\|_2^2 &\leq \frac{1}{\kappa_L^2 \lambda_{\epsilon_N}^2} \{1 + \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2 \Psi^2(\overline{\mathcal{M}}) \\ &\quad + \frac{4}{\kappa_L \lambda_{\epsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*) \\ &\leq \frac{4}{\kappa_L^2 \lambda_{\epsilon_N}^2} \Psi^2(\overline{\mathcal{M}}) + \frac{4}{\kappa_L \lambda_{\epsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (\text{A.13})$$

Case 2: $\lambda_{\epsilon_N} > 1/\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))$

Using a similar analysis as in Case 1,

$$\mathcal{R}(\hat{\Delta}_{\epsilon_N, \overline{\mathcal{M}}^\perp}) \leq \Psi(\overline{\mathcal{M}}^\perp) \|\hat{\Delta}_{\epsilon_N}\|_2. \quad (\text{A.14})$$

Substituting (A.14) and (A.10) into (A.8), we obtain,

$$\begin{aligned} 0 \geq & \kappa_L \lambda_{\epsilon_N} \|\hat{\Delta}_{\epsilon_N}\|_2^2 + \{1 - \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \Psi(\overline{\mathcal{M}}^\perp) \|\hat{\Delta}_{\epsilon_N}\|_2 \\ & - \{1 + \lambda_{\epsilon_N} \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\} \Psi(\overline{\mathcal{M}}) \|\hat{\Delta}_{\epsilon_N}\|_2 - 2\mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (\text{A.15})$$

Now using the inequality (A.12), yields,

$$\begin{aligned} \|\hat{\Delta}_{\epsilon_N}\|_2^2 &\leq \frac{1}{\kappa_L^2} \left(\left\{ \frac{1}{\lambda_{\epsilon_N}} - \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \right\} \Psi(\overline{\mathcal{M}}^\perp) \right. \\ &\quad \left. - \left\{ \frac{1}{\lambda_{\epsilon_N}} + \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \right\} \Psi(\overline{\mathcal{M}}) \right)^2 \\ &\quad + \frac{4}{\kappa_L \lambda_{\epsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned} \quad (\text{A.16})$$

Applying the inequality $(A + B)^2 \leq 2A^2 + 2B^2$ to the first term in (A.16) gives,

$$\begin{aligned} \|\hat{\Delta}_{\epsilon_N}\|_2^2 &\leq \frac{2}{\kappa_L^2} \left\{ \frac{1}{\lambda_{\epsilon_N}} - \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \right\}^2 \Psi^2(\overline{\mathcal{M}}^\perp) \\ &\quad + \frac{2}{\kappa_L^2} \left\{ \frac{1}{\lambda_{\epsilon_N}} + \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \right\}^2 \Psi^2(\overline{\mathcal{M}}) \\ &\quad + \frac{4}{\kappa_L \lambda_{\epsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned}$$

Note that $0 < 1/\lambda_{\epsilon_N} < \mathcal{R}^*(\nabla \mathcal{L}(\theta^*))$, therefore,

$$\left\{ \frac{1}{\lambda_{\epsilon_N}} - \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \right\}^2 \leq \{\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2. \quad (\text{A.17})$$

We also have,

$$\left\{ \frac{1}{\lambda_{\epsilon_N}} + \mathcal{R}^*(\nabla \mathcal{L}(\theta^*)) \right\}^2 \leq 4\{\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2. \quad (\text{A.18})$$

Therefore, combining (A.17) and (A.18),

$$\begin{aligned} \|\hat{\Delta}_{\epsilon_N}\|_2^2 &\leq \frac{2}{\kappa_L^2} \{\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2 \Psi^2(\overline{\mathcal{M}}) \\ &\quad + \frac{8}{\kappa_L^2} \{\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))\}^2 \Psi^2(\overline{\mathcal{M}}^\perp) \\ &\quad + \frac{4}{\kappa_L \lambda_{\epsilon_N}} \mathcal{R}(\theta_{\mathcal{M}^\perp}^*). \end{aligned}$$

□

A.3 Preliminary propositions for Theorem 5.1

In this Appendix we present three propositions that assist in the development of the proof of Theorem 5.1:

Proposition A.1 Consider the optimization problem in (16), and denote by λ_{ϵ_N} the corresponding Lagrange multiplier of its constraint. Then, if $\hat{\theta}_{\epsilon_N} \neq 0$, λ_{ϵ_N} can be computed as

$$\lambda_{\epsilon_N} = \frac{1}{\|\nabla \mathcal{L}(\hat{\theta}_{\epsilon_N})\|_\infty}. \quad (\text{A.19})$$

Proof. See Appendix A.4. \square

Proposition A.2 Suppose Assumptions 5.2 and 5.3 hold. Then, with probability $1 - n\beta$ ($0 \leq \beta \leq 1/n$), we have

$$P(\|\nabla \mathcal{L}(\theta^*)\|_\infty \leq t | \Phi_N) \geq \left(1 - 2 \exp \left[-\frac{N^2 t^2}{2\sigma_e^2 \chi_\beta^2(\Sigma, I)} \right] \right)^n,$$

where s_{\max} is the maximum element on the diagonal of the matrix Σ . In particular, choosing a specific value for t ,

$$\|\nabla \mathcal{L}(\theta^*)\|_\infty \leq \frac{\sqrt{2\sigma_e^2 \chi_\beta^2(\Sigma, I) \ln(2/\beta)}}{N}, \quad (\text{A.20})$$

with probability at least $1 - 2n\beta$ ($0 \leq \beta \leq 1/2n$).

Proof. See Appendix A.5. \square

Proposition A.3 Suppose Assumptions 5.2 and 5.3 hold, then with probability at least $1 - n\beta$ ($0 \leq \beta \leq 1/n$), we have

$$P\left(\left\|\nabla \mathcal{L}(\hat{\theta}_{\epsilon_N})\right\|_\infty \leq t \mid e\right) \geq \left\{1 - 2 \exp \left(-\frac{N^2 t^2}{2\sigma_e^2 s_{\max} \chi_\beta^2(N-n)(1+\epsilon_N)} \right) \right\}^n$$

where s_{\max} is the maximum element on the diagonal of the matrix Σ . In particular, choosing a specific value for t ,

$$\left\|\nabla \mathcal{L}(\hat{\theta}_{\epsilon_N})\right\|_\infty \leq \frac{\sqrt{2\sigma_e^2 s_{\max} \chi_\beta^2(N-n)(1+\epsilon_N) \ln(2/\beta)}}{N}, \quad (\text{A.21})$$

with probability at least $1 - 2n\beta$ ($0 \leq \beta \leq 1/2n$).

Proof. See Appendix A.6. \square

A.4 Proof of Proposition A.1

Let us rewrite the SPARSEVA problem (16),

$$\begin{aligned} \hat{\theta}_{\epsilon_N} \in \arg \min_{\theta \in \mathbb{R}^n} \quad & \|\theta\|_1 \\ \text{s.t.} \quad & \mathcal{L}(\theta) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N) \leq 0, \end{aligned} \quad (\text{A.22})$$

in the Lagrangian form (1) using Lemma A.2. The Lagrangian of the optimization problem (A.22) is,

$$\begin{aligned} g(\theta, \lambda) &= \|\theta\|_1 + \lambda(\mathcal{L}(\theta) - \mathcal{L}(\hat{\theta}_{NR})(1 + \epsilon_N)) \\ &= \|\theta\|_1 + \frac{\lambda}{2N}(\|Y_N - \Phi_N^T \theta\|_2^2 - \|Y_N - \Phi_N^T \hat{\theta}_{NR}\|_2^2(1 + \epsilon_N)). \end{aligned} \quad (\text{A.23})$$

The subdifferential of $g(\theta, \lambda)$ can be computed as

$$\frac{\partial g(\theta, \lambda)}{\partial \theta} = v - \frac{\lambda}{N} \Phi_N(Y_N - \Phi_N^T \theta), \quad (\text{A.24})$$

where $v = (v_1, \dots, v_m)^T$ is of the form

$$\begin{cases} v_i = 1 & \text{if } \theta_i > 0 \\ v_i = -1 & \text{if } \theta_i < 0 \\ v_i \in [-1, 1] & \text{if } \theta_i = 0. \end{cases} \quad (\text{A.25})$$

Using property 2 of Lemma A.2, when $\hat{\theta}_{\epsilon_N}$ is a solution of the SPARSEVA problem (16) and λ_{ϵ_N} is a Lagrange multiplier, we have,

$$\begin{aligned} 0 &= \frac{\partial g(\theta, \lambda)}{\partial \theta} \Big|_{\theta=\hat{\theta}_{\epsilon_N}, \lambda=\lambda_{\epsilon_N}} \\ &= -\frac{\lambda_{\epsilon_N}}{N} \Phi_N(Y_N - \Phi_N^T \hat{\theta}_{\epsilon_N}) + v_{\hat{\theta}_{\epsilon_N}}, \end{aligned} \quad (\text{A.26})$$

for some v of the form in (A.25). Note that when $\hat{\theta}_{\epsilon_N} \neq 0$, $\|v\|_\infty = 1$, which means that

$$\lambda_{\epsilon_N} = \frac{N}{\|\Phi_N(Y_N - \Phi_N^T \hat{\theta}_{\epsilon_N})\|_\infty}.$$

Since $\nabla \mathcal{L}(\hat{\theta}_{\epsilon_N}) = \frac{1}{N} \Phi_N(Y_N - \Phi_N^T \hat{\theta}_{\epsilon_N})$, we can also write λ_{ϵ_N} as

$$\lambda_{\epsilon_N} = \frac{1}{\|\nabla \mathcal{L}(\hat{\theta}_{\epsilon_N})\|_\infty}.$$

\square

Note that this proof is similar to the one in [14], where an expression was derived for the Lagrange multiplier in the traditional l_1 norm regularization problem (the LASSO). Here we have derived the Lagrange multiplier for the SPARSEVA problem as given in (16).

A.5 Proof of Proposition A.2

For the linear regression (13) and the choice of $\mathcal{L}(\theta)$ in (15),

$$\nabla \mathcal{L}(\theta^*) = \frac{1}{N} \Phi_N(Y_N - \Phi_N^T \theta^*) = \frac{1}{N} \Phi_N e.$$

Denote R_j as the j^{th} row of the matrix Φ_N , then $\nabla \mathcal{L}(\theta^*)$ can be computed as,

$$\nabla \mathcal{L}(\theta^*) = \frac{1}{N} \Phi_N e = \frac{1}{N} \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{bmatrix} e = \frac{1}{N} \begin{bmatrix} R_1 e \\ R_2 e \\ \vdots \\ R_n e \end{bmatrix}$$

consider the variable $Z = N^{-1}R_j e$, using Assumption 5.3 on the disturbance noise e , $e \sim \mathcal{N}(0, \sigma_e^2)$, we have,

$$Z|e \sim \mathcal{N}\left(0, \frac{\sigma_e^2}{N^2}R_j R_j^T\right). \quad (\text{A.27})$$

Now in order to derive a bound for $\nabla \mathcal{L}(\theta^*)$, we first derive an upper bound for the variance $N^{-2}\sigma_e^2 R_j R_j^T$ of the distribution in (A.27). Since $R_j \sim \mathcal{N}(0, \Sigma)$, we have,

$$R_j R_j^T \sim \chi^2(\Sigma, I),$$

where $\chi^2(\Sigma, I)$ is the generalized Chi squared with parameters Σ and I . Hence, with probability $1 - \beta$, $0 \leq \beta \leq 1$, we have,

$$R_j R_j^T \leq \chi_\beta^2(\Sigma, I). \quad (\text{A.28})$$

Hence, the variance of the distribution of the variable $N^{-1}R_j e$, is,

$$\frac{\sigma_e^2}{N^2}R_j R_j^T \leq \frac{\sigma_e^2}{N^2}\chi_\beta^2(\Sigma, I), \quad (\text{A.29})$$

with probability $1 - \beta$.

Note that from (A.27), for any $t > 0$, we have,

$$P\left(\left|\frac{1}{N}R_j e\right| \leq t \mid \Phi_N\right) = \int_{-t}^t f\left(x \mid 0, \frac{\sigma_e^2}{N^2}R_j R_j^T\right) dx, \quad (\text{A.30})$$

where $f(x|0, N^{-2}\sigma_e^2 R_j R_j^T)$ denotes the pdf of the Normal distribution $\mathcal{N}(0, N^{-2}\sigma_e^2 R_j R_j^T)$. This gives,

$$P\left(\left\|\frac{\Phi_N e}{N}\right\|_\infty \leq t \mid \Phi_N\right) = \prod_{j=1}^n \left\{ \int_{-t}^t f\left(x \mid 0, \frac{\sigma_e^2}{N^2}R_j R_j^T\right) dx \right\}. \quad (\text{A.31})$$

This expression can be bounded from below using the standard result that $P(|\mathcal{N}(0, \sigma^2)| > t) \leq 2\exp(-t^2/2\sigma^2)$ [21, Eq. (5.5)], to obtain

$$P\left(\left\|\frac{\Phi_N e}{N}\right\|_\infty \leq t \mid \Phi_N\right) \geq \prod_{j=1}^n \left(1 - 2\exp\left[-\frac{N^2 t^2}{2\sigma_e^2 R_j R_j^T}\right]\right).$$

The expression in parentheses on the right hand side is monotonically decreasing in $R_j R_j^T$, so using (A.29) gives

$$P\left(\left\|\frac{\Phi_N e}{N}\right\|_\infty \leq t \mid \Phi_N\right) \geq \left(1 - 2\exp\left[-\frac{N^2 t^2}{2\sigma_e^2 \chi_\beta^2(\Sigma, I)}\right]\right)^n,$$

which holds with probability² at least $1 - n\beta$.

² This bound follows because the events A_j that (A.29) holds are not necessarily independent, but their joint probability can be bounded like $P(A_1 \cap \dots \cap A_n) = 1 - P(A_1^C \cup \dots \cup A_n^C) \geq 1 - P(A_1^C) - \dots - P(A_n^C) = 1 - n\beta$.

In particular, taking $t = \sqrt{2\sigma_e^2 \chi_\beta^2(\Sigma, I) \ln(2/\beta)}/N$ gives

$$P\left(\left\|\frac{\Phi_N e}{N}\right\|_\infty \leq \frac{\sqrt{2\sigma_e^2 \chi_\beta^2(\Sigma, I) \ln(2/\beta)}}{N} \mid \Phi_N\right) \geq (1 - \beta)^n \geq 1 - n\beta$$

with probability at least $1 - n\beta$, or equivalently,

$$\left\|\nabla \mathcal{L}(\theta^*)\right\|_\infty \leq \frac{\sqrt{2\sigma_e^2 \chi_\beta^2(\Sigma, I) \ln(2/\beta)}}{N}$$

with probability at least $1 - 2n\beta$. \square

A.6 Proof of Proposition A.3

When $\hat{\theta}_{\varepsilon_N}$ is the solution of the problem in (16), we have,

$$\nabla \mathcal{L}(\hat{\theta}_{\varepsilon_N}) = \frac{1}{N}\Phi_N(Y_N - \Phi_N^T \hat{\theta}_{\varepsilon_N}). \quad (\text{A.32})$$

Denote $e_{\varepsilon_N} = Y_N - \Phi_N^T \hat{\theta}_{\varepsilon_N}$, and R_j as the j^{th} row of the matrix Φ_N , then (A.32) becomes,

$$\nabla \mathcal{L}(\hat{\theta}_{\varepsilon_N}) = \frac{1}{N}\Phi_N e_{\varepsilon_N} = \frac{1}{N} \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{bmatrix} e_{\varepsilon_N} = \frac{1}{N} \begin{bmatrix} R_1 e_{\varepsilon_N} \\ R_2 e_{\varepsilon_N} \\ \vdots \\ R_n e_{\varepsilon_N} \end{bmatrix}.$$

From Assumption 5.2, and using the same argument as in Proposition A.2, we have that each element of R_j is distributed as $\mathcal{N}(0, \Sigma(j, j))$.

Consider the variable $Z = N^{-1}e_{\varepsilon_N}^T R_j^T$, Since $R_j \sim \mathcal{N}(0, \Sigma)$,

$$Z \sim \mathcal{N}\left(0, \frac{1}{N^2}e_{\varepsilon_N}^T \Sigma e_{\varepsilon_N}\right).$$

Since Σ is symmetric and positive definite matrix, hence using singular value decomposition, we can find a diagonal matrix D that satisfies,

$$\Sigma = Q^T D Q, \quad (\text{A.33})$$

where Q is the unitary matrix, i.e. $Q Q^T = I$. Therefore, we have,

$$\frac{1}{N^2}e_{\varepsilon_N}^T \Sigma e_{\varepsilon_N} \leq \frac{s_{\max}}{N^2}e_{\varepsilon_N}^T e_{\varepsilon_N}, \quad (\text{A.34})$$

where s_{\max} is the maximum element on the diagonal of matrix D , i.e. maximum singular value of matrix Σ . Note

that,

$$\begin{aligned} e_{\varepsilon_N}^T e_{\varepsilon_N} &= (Y_N - \Phi_N^T \hat{\theta}_{\varepsilon_N})^T (Y_N - \Phi_N^T \hat{\theta}_{\varepsilon_N}) = 2N\mathcal{L}(\hat{\theta}_{\varepsilon_N}) \\ &= 2N\mathcal{L}(\hat{\theta}_{NR})(1 + \varepsilon_N). \end{aligned} \quad (\text{A.35})$$

From Section 4.4 in [18], we have,

$$\mathcal{L}(\hat{\theta}_{NR})|\Phi_N \sim \frac{\sigma_e^2}{2N}\chi^2(N-n),$$

which gives,

$$\mathcal{L}(\hat{\theta}_{NR}) \leq \frac{\sigma_e^2}{2N}\chi_{\beta}^2(N-n),$$

with probability $1 - \beta$, $0 \leq \beta \leq 1$. Combining this inequality with (A.34) and (A.35) gives, with probability $1 - \beta$,

$$\frac{1}{N^2} e_{\varepsilon_N}^T \Sigma e_{\varepsilon_N} \leq \frac{\sigma_e^2}{N^2} s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N).$$

Hence,

$$\begin{aligned} P\left(\left|\frac{R_j e_{\varepsilon_N}}{N}\right| \leq t \mid e\right) \\ \geq \int_{-t}^t f\left(x \mid 0, s_{\max} \frac{\sigma_e^2}{N^2} \chi_{\beta}^2(N-n)(1 + \varepsilon_N)\right) dx, \end{aligned} \quad (\text{A.36})$$

with probability $1 - \beta$. This means,

$$\begin{aligned} P\left(\left\|\frac{\Phi_N e_{\varepsilon_N}}{N}\right\|_{\infty} \leq t \mid e\right) \\ = \prod_{j=1}^n \left\{ \int_{-t}^t f\left(x \mid 0, \frac{\sigma_e^2}{N^2} s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N)\right) dx \right\} \\ \geq \left\{ 1 - 2 \exp\left(-\frac{N^2 t^2}{2\sigma_e^2 s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N)}\right) \right\}^n \end{aligned} \quad (\text{A.37})$$

with probability at least $1 - n\beta$, following the same reasoning as in the proof of Proposition A.2. Therefore,

$$\begin{aligned} P\left(\left\|\nabla \mathcal{L}(\hat{\theta}_{\varepsilon_N})\right\|_{\infty} \leq t \mid e\right) \\ \geq \left\{ 1 - 2 \exp\left(-\frac{N^2 t^2}{2\sigma_e^2 s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N)}\right) \right\}^n \end{aligned}$$

with probability at least $1 - n\beta$.

Taking $t = \sqrt{2\sigma_e^2 s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N) \ln(2/\beta)}/N$ gives

$$\begin{aligned} P\left(\left\|\nabla \mathcal{L}(\hat{\theta}_{\varepsilon_N})\right\|_{\infty} \leq \frac{\sqrt{2\sigma_e^2 s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N) \ln(2/\beta)}}{N} \mid e\right) \\ \geq (1 - \beta)^n \geq 1 - n\beta, \end{aligned}$$

with probability at least $1 - n\beta$, or, equivalently,

$$\left\|\nabla \mathcal{L}(\hat{\theta}_{\varepsilon_N})\right\|_{\infty} \leq \frac{\sqrt{2\sigma_e^2 s_{\max} \chi_{\beta}^2(N-n)(1 + \varepsilon_N) \ln(2/\beta)}}{N},$$

with probability at least $1 - 2n\beta$. \square