

1 **SPARSE PARAMETER IDENTIFICATION FOR STOCHASTIC**
2 **SYSTEMS BASED ON L_γ REGULARIZATION***

3 JIAN GUO[†], YING WANG[†], YANLONG ZHAO[†], AND JI-FENG ZHANG^{†‡}

4 **Abstract.** This paper is concerned with the reconstruction of the zero and non-zero elements
5 of the sparse parameter vector of stochastic systems with general observation sequences. A sparse
6 parameter identification algorithm based on L_γ penalty with $0 < \gamma < 1$ and the residual sum of
7 squares is proposed. Without requiring independently and identically distributed (i.i.d) and station-
8 ary conditions on the observation sequences, the proposed algorithm is proved that not only the
9 contributing variable corresponding to the non-zero parameters can be selected out with probability
10 converging to one, but also the estimates of the non-zero parameters have the asymptotic normality
11 property. In order to improve the performance of the L_γ regularization method, a two-step algorithm
12 based on adaptively weighted L_γ penalty with $0 < \gamma \leq 1$ is designed, whose set and parameter al-
13 most sure convergence are established with non-i.i.d and non-stationary observation sequences. The
14 proposed methods are applied to the structure selection of the nonlinear autoregressive models with
15 exogenous variables and the sparse parameter identification of the linear feedback control systems.
16 Finally, three numerical examples are given to verify the efficiency of the theoretical results.

17 **Key word.** Stochastic system, sparse identification, L_γ penalty, asymptotic normality, strong
18 consistency.

19 **MSC codes.** 93E03, 93E12, 93E24

20 **1. Introduction.** The sparsity problems are occurring in many areas of scien-
21 tific research and engineering practice and have attracted considerable attention in
22 recent years. Exemplary applications involve image processing [25, 30], wireless com-
23 munication [13, 27], biometrics [31], compressed sampling [3], and so on. One of the
24 most interesting issues is the exact reconstruction of zero and non-zero elements of
25 sparse parameter vectors. This is of great importance in engineering applications as it
26 provides a way to implement a parsimonious model with better predictive performance
27 and can reduce the curse of dimensionality.

28 Classical parameter identification is a rapidly developing field for the reconstruc-
29 tion of system elements and has achieved a great success in both theoretical research
30 and practical applications [6, 26]. A series of prestigious methods have been devel-
31 oped, including stochastic gradient descent, stochastic approximation, least-squares
32 (LS), least mean square, and so on. These methods are usually obtained by min-
33 imizing some criteria such as the square error between the predicted and observed
34 signals, and have some theoretical properties, such as consistency, convergence rate,
35 asymptotic normality, etc. However, for sparse systems, since they tend to be high-
36 dimensional or have a limited number of samples, these classical theories and methods
37 will no longer be valid.

38 In the field of statistics, a number of effective and widely used methods have
39 emerged for sparse problems [9]. For instance, there are several classical criteria to
40 implement variable selection, such as Akaike’s information criterion (AIC) [1] and

*Submitted to the editors on September 7, 2023.

Funding: The work is supported by National Key R&D Program of China under Grant 2018YFA0703800, National Natural Science Foundation of China under Grant T2293770, 62025306, 62303452 and 12226305, CAS Project for Young Scientists in Basic Research under Grant YSBR-008 and China Postdoctoral Science Foundation under Grant 2022M720159.

[†]Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, China and School of Mathematics Sciences, University of Chinese Academy of Sciences, Beijing 100149, China (j.guo@amss.ac.cn, wangying96@amss.ac.cn, ylzhao@amss.ac.cn, jif@iss.ac.cn).

[‡]Corresponding author.

41 Bayesian information criterion (BIC) [32]. However, they are not applicable to high-
 42 dimensional data as they may involve solving NP-Hard optimization problems. Subse-
 43 quently, regularization methods are proposed and widely used as a solution to sparse
 44 problem. Typically, regularization is designed by adding a penalty term to the LS
 45 objective, where the penalty term is generally defined as a norm over the parameter
 46 space. L_0 regularization is the first regularization method applied to variable selec-
 47 tion, which can produce the sparsest solution, but requires solving a combinatorial
 48 optimization problem, whose complexity grows exponentially with dimension. [33]
 49 proposed an alternative method called LASSO which converts the combinatorial opti-
 50 mization problem of variable selection into an easily solvable quadratic programming
 51 problem, but is not as sparse as L_0 regularization. Thereafter, various regularization
 52 methods such as smoothly clipped absolute deviation [9], adaptive LASSO [43], elastic
 53 net [44], etc., have become the main tools for data analysis. In addition, [18] consid-
 54 ers the asymptotic behavior of regression estimates that minimize the sum of squared
 55 residuals plus the L_γ penalty. [37] and [38] addressed the particular importance of
 56 $L_{1/2}$ regularization in sparse modeling and obtained promising practical results in
 57 image processing, matrix filling, etc.

58 With the rapid development of variable selection in statistics, some of these ideas
 59 and methods have been applied to stochastic systems and control. For instance, [34]
 60 used the L_0 regularization to obtain the sparsest estimate of the parameter vector.
 61 [24] utilized L_1 regularization to identify the system parameters and predict future
 62 signals assuming that the output noise components exhibited strong seriality and
 63 cross-sectional correlation. [42] introduced an LS sparse parameter identification al-
 64 gorithm based on L_1 penalty with adaptive weights and proved its convergence with
 65 general observation sequences, and then [12] generalized this approach to Multivariate
 66 ARMA Systems with Exogenous Inputs. In addition, some non-convex regularization
 67 methods are also employed for stochastic systems. [11] suggested a simple numerical
 68 scheme to compute solutions with minimal L_γ norm and studied its convergence. [29]
 69 proposed a new sparse signal reconstruction algorithm based on the minimization of
 70 the squared error of a smooth L_γ ($\gamma < 1$) norm regularization, which provided bet-
 71 ter signal reconstruction performance. [36] presented generalized shrinkage penalties
 72 with explicit proximal mappings and thus gave iterative γ -shrinkage iterative algo-
 73 rithms that could be implemented to accurately recover a given sparse data with a
 74 given measurement matrix. However, these papers about non-convex penalized meth-
 75 ods, do not give theoretical results like that in [42]: whether the solutions obtained
 76 by non-convex regularization methods are still convergent in the non-stationary and
 77 non-independently and identically distributed (i.i.d) situation.

78 **Motivation of this work.** As known in the literature, the L_1 regularization
 79 method has led to remarkable progress in sparse problems. However, L_1 regulariza-
 80 tion suffers from bias, leading to a heavily biased estimate and not achieving reliable
 81 recovery with the least observations [4]. Besides, L_1 regularization may produce in-
 82 consistent selections when applied to some situations [43]. In contrast, the non-convex
 83 penalty such as L_γ ($0 < \gamma < 1$) regularization has the advantage of improving the bias
 84 problem and has led to significant performance improvements in many applications.
 85 For instance, [19] demonstrated the very high efficiency of applying $L_{1/2}$ and $L_{2/3}$
 86 regularization to image deconvolution. This motivates us to investigate non-convex
 87 penalties in the fields of systems and control. However, in the existing literature on
 88 the non-convex regularization, the noise is usually required to be i.i.d or there is prior
 89 knowledge of the sample probability distribution, or the observed sequences are deter-
 90 ministic [10]. These conditions are difficult to satisfy for stochastic systems, especially

91 feedback control systems. Besides, it is not clear whether the estimates obtained by
 92 utilizing such an approach in sparse system identification still have the theoretical
 93 asymptotic properties.

94 Thus, this paper sets out to investigate the non-convex $L_\gamma(0 < \gamma < 1)$ regular-
 95 ization method in sparse identification problems of stochastic dynamic systems with
 96 general observation sequences and non-i.i.d noise. The main contributions of this
 97 paper are as follows:

- 98 • This paper proposes a sparse parameter identification algorithm based on the L_γ
 99 ($0 < \gamma < 1$) penalty and the residual sum of squares for stochastic sparse systems
 100 with non-i.i.d and non-stationary observation sequences and non-i.i.d noise. This
 101 algorithm yields significantly better performance in terms of sparsity induction and
 102 efficiency compared to the convex penalty. In addition, the theoretical properties
 103 of this algorithm are established. Specifically, the almost sure convergence of the
 104 estimates is proven. Besides, the set convergence in probability is shown, i.e., the
 105 probability that the proposed algorithm correctly selects the non-zero elements of
 106 the unknown sparse parameter vector converges to one. Moreover, the asymptotic
 107 normality of the parameter estimates is obtained. These results incorporate the
 108 results of bridge estimate [18] and do not require additional strong irrepresentable
 109 conditions compared with LASSO [40].
- 110 • In order to improve the performance of the L_γ regularization method, motivated by
 111 [42] and [43], a two-step algorithm based on the adaptively weighted $L_\gamma(0 < \gamma \leq 1)$
 112 penalty and the residual sum of squares is proposed. For the case of non-i.i.d
 113 and non-stationary observation sequences and non- i.i.d noise, not only is almost
 114 sure parameter convergence established, but also almost sure set convergence is
 115 achieved, i.e., this algorithm correctly selects the non-zero elements of the unknown
 116 sparse parameter vector with probability one using a finite number of observations.
 117 Moreover, this algorithm is more efficient in sparsity induction than the adaptive
 118 LASSO and the algorithm in [42] and covers their results when $\gamma = 1$.
- 119 • The proposed sparse identification algorithms in this paper are applied to two kinds
 120 of typical scenes in stochastic sparse systems with non-i.i.d observation sequences.
 121 Specifically, the proposed algorithms can efficiently select the contributing basis
 122 functions out for the Nonlinear AutoRegressive models with eXogenous variables
 123 (NARX). Furthermore, the proposed algorithm is able to accurately reconstruct
 124 the sparse parameters of the linear feedback control systems with non-i.i.d and
 125 non-stationary observation sequences and non-i.i.d noise.

126 The rest of this paper is organized as follows: Section 2 gives the problem formu-
 127 lation. Section 3 proposes the $L_\gamma(0 < \gamma < 1)$ regularization algorithm, establishes its
 128 theoretical results and compares it with related works. Section 4 gives an adaptively
 129 weighted two-step algorithm and investigates its properties. In Section 5, the pro-
 130 posed algorithm is applied to accomplish the structure selection of the NARX model
 131 and the sparse identification of the linear feedback control systems. In Section 6,
 132 three typical simulation examples are given to illustrate the algorithms' performance.
 133 And in Section 7, some concluding remarks and further works are provided.

134 **Notation:** Let $(\Omega, \mathcal{F}, \mathbb{P})$ be the probability space, $\omega \in \Omega$ be the sample points,
 135 and $E(\cdot)$ be the expectation operator. $\|\cdot\|_1$ and $\|\cdot\|$ denote 1-norm and 2-norm for
 136 vectors or matrices, respectively. By \mathbb{R} and \mathbb{N}_+ , we denote the sets of real numbers
 137 and positive integers, respectively. \mathbf{I}_p represents a unit matrix of order p and $\mathbf{0}_p =$
 138 $[0, \dots, 0]^T \in \mathbb{R}^p$. Moreover, $\text{sign}(\cdot)$ is defined as $\text{sign}(x) = 1$, when $x \geq 0$, and $\text{sign}(x) =$
 139 -1 , when $x < 0$, $\text{vec}(x_j)_{j=1}^q$ means $[x_1, x_2, \dots, x_q]^T$, and for a set A , by A^c , we
 140 denote the complement of A . For any two positive sequences $\{a_k\}_{k \geq 1}$ and $\{b_k\}_{k \geq 1}$,

141 $a_k = O(b_k)$ means there are $c > 0$ and $k_0 \in \mathbb{N}_+$ such that $a_k \leq cb_k$ for all $k \geq k_0$;
 142 $a_k = o(b_k)$ means $a_k/b_k \rightarrow 0$ as $k \rightarrow \infty$. For two random sequences $\{x_k\}$ and $\{y_k\}$,
 143 we give the following two frequently-used definitions in this paper

- 144 • $x_k = O_p(y_k)$ means that for any $\epsilon > 0$, there is a finite $M > 0$ and a finite $N > 0$
 145 such that $\mathbb{P}\{|x_k| \geq M|y_k|\} < \epsilon$ for all $k \geq N$;
- 146 • $x_k = o_p(y_k)$ means $x_k/y_k \xrightarrow{P} 0$ as $k \rightarrow \infty$, where \xrightarrow{P} means convergence in
 147 probability.

148 **2. Problem formulation.** Consider the stochastic sparse system

$$149 \quad (2.1) \quad y_{k+1} = \theta^T \varphi_k + w_{k+1}, \quad k \geq 0,$$

150 where $\theta = [\theta(1), \dots, \theta(p)]^T \in \mathbb{R}^p$ is the unknown p -dimensional parameter vector
 151 containing many zero values, $\varphi_k \in \mathbb{R}^p$ consisting of possibly current and past inputs
 152 and outputs, is the stochastic regressor vector, y_{k+1} and w_{k+1} are the system output
 153 and noise, respectively. Denote the zero elements set of the unknown parameter θ by
 154 $A^* = \{j : \theta(j) = 0, j \in \{1, \dots, p\}\}$. Suppose that there are q non-zero elements in
 155 the vector θ . Without loss of generality, we assume that $\theta(j) = 0$ for $j = q+1, \dots, p$.

156 **Problem.** The identification problem of the stochastic sparse system is to infer
 157 the zero elements A^* and to estimate the non-zero elements of the unknown parameter
 158 vector θ by using the observed data $\{y_{k+1}, \varphi_k\}_{k=1}^n$.

159 Before giving the sparse identification algorithm, the following assumptions are
 160 introduced.

Assumptions. Denote the family of the σ -algebras $\{\mathcal{F}_k\}$ as

$$\mathcal{F}_k = \sigma\{\varphi_0, \dots, \varphi_k, w_1, \dots, w_k\}, \quad k \geq 1,$$

161 the maximum and minimum eigenvalues of $\sum_{k=1}^n \varphi_k \varphi_k^T$ as $\lambda_{\max}(n)$ and $\lambda_{\min}(n)$, re-
 162 spectively, and the maximum eigenvalue of $E \sum_{k=1}^n \varphi_k \varphi_k^T$ as $\lambda_{E, \max}(n)$.

163 **(A1)** The noise $\{w_k, \mathcal{F}_k\}_{k \geq 1}$ is a martingale difference sequence and there is $\delta > 0$

164 such that $\sup_k E[|w_{k+1}|^{2+\delta} | \mathcal{F}_k] < \infty$, a.s.

165 **(A2)** (a) For the maximal and minimal eigenvalues of $\sum_{k=1}^n \varphi_k \varphi_k^T$, it holds

$$\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} \xrightarrow{n \rightarrow \infty} 0 \text{ a.s.}$$

166 (b) For each n , there is a positive number d_n such that

$$d_n \lambda_{\min}(n)^{-1} = O_p(1) \text{ and } \frac{\sqrt{\lambda_{E, \max}(n)}}{d_n} \xrightarrow{n \rightarrow \infty} 0.$$

167 *Remark 2.1.* In Assumption (A1), a sequence of martingale differences is broader
 168 than a sequence of independent variables, which implies a much milder restriction
 169 on sequence memory than independence and allows w_{k+1} to depend on \mathcal{F}_k . Many
 170 random variables, such as Gaussian random variables, uniformly distributed random
 171 variables, and so on, all satisfy this assumption. Assumptions (A2) is about the
 172 system observation sequences. Assumption (A2)(a) is the classical weakest strong
 173 convergence condition for LS [22].

174 **3. L_γ regularization algorithm and its properties.** This section constructs
 175 a sparse identification algorithm based on L_γ ($0 < \gamma < 1$) regularization and gives the
 176 corresponding theoretical properties.

177 **3.1. L_γ regularization algorithm.** We start by giving the objective function
 178 based on L_γ penalty with $0 < \gamma < 1$ and residual sum of squares:

$$179 \quad (3.1) \quad J_n(\beta) = \sum_{k=1}^n (y_{k+1} - \beta^T \varphi_k)^2 + \lambda_n \sum_{l=1}^p |\beta(l)|^\gamma,$$

180 where λ_n is a penalty parameter and $\beta = [\beta(1), \dots, \beta(p)]^T$.

181 **Assumption.** We first give the following assumption about the parameter λ_n .

182 **(A3)** The penalty parameter $\{\lambda_n\}$ of (3.1) satisfies that

$$183 \quad (a) \frac{\lambda_n}{\lambda_{\min}(n)} \xrightarrow{n \rightarrow \infty} 0, \text{ a.s.}, \quad (b) \frac{\lambda_n}{\lambda_{E, \max}(n)^{1/2}} \xrightarrow{n \rightarrow \infty} 0, \quad (c) \frac{\lambda_n d_n^{2-\gamma}}{\lambda_{E, \max}(n)^{2-\frac{1}{2}\gamma}} \xrightarrow{n \rightarrow \infty} \infty.$$

184 *Remark 3.1.* Assumption (A3) is about the penalty parameter λ_n . It is able to
 185 be satisfied and cover the classical persistent excitation condition as a special case,
 186 i.e., $C_1 n \leq \lambda_{\min}(n) \leq \lambda_{\max}(n) \leq C_2 n$ for some constants C_1 and C_2 . Specifically,
 187 d_n in (A2)(b) can be n and for any given $0 < \gamma < 1$, λ_n can be chosen as n^α with
 188 $\frac{1}{2}\gamma < \alpha < \frac{1}{2}$ to meet Assumption (A3).

189 **Algorithm.** The sparse identification algorithm based on L_γ penalty is designed
 190 in Algorithm 3.1. This algorithm provides a method for combining variable selection
 and parameter estimation in a single step.

Algorithm 3.1 L_γ regularization.

Step 0 (Initialization). For given $0 < \gamma < 1$, choose a positive sequence $\{\lambda_n\}_{n \geq 1}$
 satisfying (A3).

Step 1 (Sparse Optimization with L_γ penalty) With γ and λ_n , optimize the
 objective function

$$(3.2) \quad J_n(\beta) = \sum_{k=1}^n (y_{k+1} - \beta^T \varphi_k)^2 + \lambda_n \sum_{l=1}^p |\beta(l)|^\gamma,$$

and obtain

$$(3.3) \quad \beta_n = [\beta_n(1), \dots, \beta_n(p)]^T = \underset{\beta}{\operatorname{argmin}} J_n(\beta),$$

$$(3.4) \quad A_n^* = \{j : \beta_n(j) = 0, j \in \{1, \dots, p\}\}.$$

191

192 *Remark 3.2.* We now discuss the feasibility of (3.3). First, the global minimum
 193 point of non-convex function $J_n(\beta)$ exists (not infinity). This is because $J_n(\beta)$ is
 194 continuous, there exists a minimum point on any compact set; and since $\|\beta\| \rightarrow \infty$,
 195 $J_n(\beta) \rightarrow \infty$, the point that minimizes $J_n(\beta)$ must be finite. Thus, (3.3) is a well-
 196 defined estimator. Second, we present the computation methods of (3.3). It is worth
 197 noting that the standard gradient-based method fails to solve this problem, because
 198 the penalty objective function $J_n(\beta)$ is non-differentiable when β has zero compo-
 199 nents. While, a large number of approximate algorithms and nonconvex optimization
 200 solvers have emerged to solve this problem. For instance, [37] proposed an iterative
 201 half thresholding algorithm for fast solution of $L_{1/2}$ regularization, and [20] and [29]
 202 designed solving algorithms by approximating the L_γ penalty with a function that
 203 has finite gradient at zero. In addition, genetic algorithms, particle swarm algorithms,
 204 simulated annealing algorithms, etc. can be used to solve non-convex optimization

205 problems as well as solvers such as IPOPT [35]. However, none of the above methods
 206 provide sufficient theoretical support. Thus, the focus of this paper is not on the
 207 discussion of the solution method of (3.3), but on the properties of its solution.

208 *Remark 3.3.* The currently existing papers on the sparse identification of L_γ pen-
 209 alty either lack theory, as in the papers [11, 29, 36], or discuss its properties only under
 210 the i.i.d and stationary condition, as in the papers [9, 18, 38]. However, in the fields of
 211 system and control, the regressor φ_k is generally non-stationary and non-independent
 212 because any real feedback controller depends essentially on the system output and
 213 hence the driven noise [17]. The main point of interest in this paper is whether the
 214 estimates (3.3) and (3.4) remain parameter convergence, set convergence and asymp-
 215 totically normality under non-stationary and non-independent conditions.

216 *Remark 3.4.* [42] proved the convergence of L_1 penalty with adaptive weights
 217 under non-stationary and non-independent assumption. While, L_γ penalty is more
 218 efficient in sparsity induction than L_1 penalty. We give an example to explain. Con-
 219 sider the Auto Regression with eXtra input (ARX) system: $y_{k+1} = \theta_1 y_k + \theta_2 u_k + w_{k+1}$
 220 with the true parameters $\theta_1 = 1$ and $\theta_2 = 0$. Let $\beta = [\beta_1, \beta_2]^T$. By Lagrange's multi-
 221 plier method, the regularized LS problem (3.3) is equivalent to solving:

$$222 \min_{\beta} J(\beta) = \sum_{k=1}^n (y_{k+1} - \beta_1 y_k - \beta_2 u_k)^2 \quad \text{s.t.} \quad |\beta_1|^\gamma + |\beta_2|^\gamma \leq s,$$

223 for some $s > 0$. Fig. 1 shows the objective function equivalence graphs of L_1 and
 224 $L_{1/2}$ penalties. The constraint region of the L_1 penalty is a square after rotation,
 225 and the constraint region of the $L_{1/2}$ penalty is a graph concave inward. The solution to
 226 this problem occurs when the contour $J(\beta)$ is first tangent to the constraint region.
 227 It can be seen that the solutions of both L_1 and $L_{1/2}$ penalties may appear at the
 228 corners, which leads to a sparse solution. This geometrically demonstrates the spar-
 229 sity of $L_\gamma (0 < \gamma \leq 1)$ regularization. Moreover, the solution of the $L_{1/2}$ regularized
 230 LS problem is more likely to appear at the corners, which implies that the solution of
 the $L_{1/2}$ regularized LS problem is sparser than L_1 .

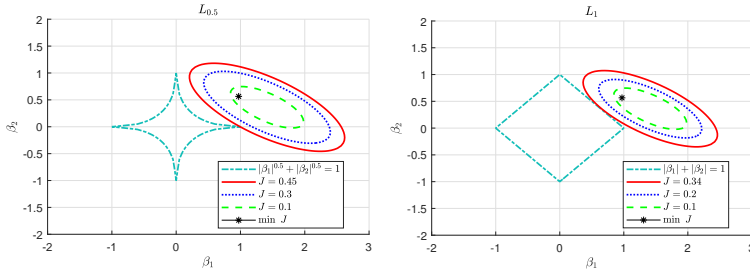


FIG. 1. L_1 penalty v.s. $L_{1/2}$ penalty

231

232 *Remark 3.5.* Now we give a way of choosing penalty parameter λ_n in Algorithm
 233 3.1 for general cases. If $\frac{\lambda_{E,\max}(n)^{3/2-\gamma/2}}{d_n^{2-\gamma}} \xrightarrow{n \rightarrow \infty} 0$ and $\frac{\sqrt{\lambda_{E,\max}(n)}}{\lambda_{\min}(n)} = O(1)$ a.s., then,
 234 for any $0 < \beta < 1$, $\lambda_n = \frac{\lambda_{E,\max}(n)^{(\frac{3}{2} - \frac{1}{2}\gamma)\beta + 1/2}}{d_n^{(2-\gamma)\beta}}$ satisfies Assumption (A3).

235 **3.2. Theoretical properties.** This section will give the theoretical properties
 236 of Algorithm 3.1. To prove these properties, we first give the following proposition.

237 PROPOSITION 3.6. [23] For the system (2.1), if Assumptions (A1) and (A2) hold,
 238 then $\left\| \left(\sum_{k=1}^n \varphi_k \varphi_k^T \right)^{-\frac{1}{2}} \sum_{k=1}^n \varphi_k w_{k+1} \right\| = O \left(\sqrt{\log \lambda_{\max}(n)} \right)$, a.s.

239 For the estimate β_n and A_n^* generated by Algorithm 3.1, the following theorem
 240 shows the almost sure convergence of the estimates.

241 THEOREM 3.7. Under Assumptions (A1), (A2)(a) and (A3)(a), the estimate given
 242 by Algorithm 3.1 is almost surely convergent, i.e., $\lim_{n \rightarrow \infty} \beta_n = \theta$, a.s.

243 *Proof.* Noting that β_n is the minimizer of $J_n(\beta)$ in Algorithm 3.1, we have
 244 $J_n(\beta_n) \leq J_n(\theta)$. Since $\lambda_n > 0$, $|\beta_n(j)|^\gamma \geq 0$, by (3.2) and a direct calculation,
 245 we have

$$\begin{aligned} 246 \quad & \lambda_n \sum_{j=1}^p |\theta(j)|^\gamma \geq \sum_{i=1}^n (y_{i+1} - \varphi_i^T \beta_n)^2 - \sum_{i=1}^n (y_{i+1} - \varphi_i^T \theta)^2 \\ 247 \quad (3.5) \quad & = (\beta_n - \theta)^T \sum_{i=1}^n (\varphi_i \varphi_i^T) (\beta_n - \theta) + 2 \sum_{i=1}^n \varphi_i^T (\theta - \beta_n) w_{i+1}. \end{aligned}$$

248 Let $P_n = \sum_{i=1}^n \varphi_i \varphi_i^T$, $\delta_n = P_n^{1/2} (\beta_n - \theta)$ and $Q_n = \left(\sum_{k=1}^n \varphi_k \varphi_k^T \right)^{-\frac{1}{2}} \sum_{k=1}^n \varphi_k w_{k+1}$.
 249 Then, (3.5) becomes

$$\begin{aligned} 250 \quad & (\beta_n - \theta)^T \sum_{i=1}^n (\varphi_i \varphi_i^T) (\beta_n - \theta) + 2 \sum_{i=1}^n \varphi_i^T (\theta - \beta_n) w_{i+1} \\ 251 \quad (3.6) \quad & = \delta_n^T \delta_n - 2 \left[\left(\sum_{i=1}^n \varphi_i \varphi_i^T \right)^{-1/2} \sum_{i=1}^n \varphi_i w_{i+1} \right]^T \delta_n = \delta_n^T \delta_n - 2 Q_n^T \delta_n. \end{aligned}$$

252 From (3.5) and (3.6) it follows that $\delta_n^T \delta_n - 2 Q_n^T \delta_n - \lambda_n \sum_{j=1}^p |\theta(j)|^\gamma \leq 0$, which implies
 253 $\|\delta_n - Q_n\|^2 - \|Q_n\|^2 - \lambda_n \sum_{j=1}^p |\theta(j)|^\gamma \leq 0$. Hence, we have

$$\begin{aligned} & \|\delta_n - Q_n\| \leq \sqrt{\|Q_n\|^2 + \lambda_n \sum_{j=1}^p |\theta(j)|^\gamma} \\ & \leq \sqrt{\|Q_n\|^2 + \lambda_n \sum_{j=1}^p |\theta(j)|^\gamma + 2 \|Q_n\| \left(\lambda_n \sum_{j=1}^p |\theta(j)|^\gamma \right)^{1/2}} = \|Q_n\| + \left(\lambda_n \sum_{j=1}^p |\theta(j)|^\gamma \right)^{1/2}. \end{aligned}$$

254 Then, by the triangular inequality we have $\|\delta_n\| \leq \|\delta_n - Q_n\| + \|Q_n\| \leq 2\|Q_n\| +$
 255 $\left(\lambda_n \sum_{j=1}^p |\theta(j)|^\gamma \right)^{1/2}$. Noting Proposition 3.6 and $\lambda_n \sum_{j=1}^p |\theta(j)|^\gamma = O(\lambda_n)$, it follows
 256 that $\|\beta_n - \theta\| = O \left(\sqrt{\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)}} + \sqrt{\frac{\lambda_n}{\lambda_{\min}(n)}} \right)$ a.s. By Assumptions (A2)(a) and
 257 (A3)(a), the proof is completed. \square

258 Next, we discuss the set convergence in probability of the estimates, starting
 259 with the following lemma to illustrate the convergence properties in probability of the
 260 estimation error.

261 LEMMA 3.8. If Assumptions (A1), (A2) and (A3)(a)-(b) hold, then

$$262 \quad (3.7) \quad \|\beta_n - \theta\| = O_p \left(\frac{\sqrt{q \lambda_{E, \max}(n)}}{d_n} \right).$$

263 *Proof.* To prove $\|\beta_n - \theta\| = O_p(\sqrt{q\lambda_{E,\max}(n)}/d_n)$, it is sufficient to prove that for
 264 any $\epsilon_1 > 0$, there exists a finite $\tilde{M} > 0$ and N such that for any $n > N$, $P(\|\beta_n - \theta\| >$
 265 $\tilde{M}\sqrt{q\lambda_{E,\max}(n)}/d_n) \leq \epsilon_1$. Let $h_n = d_n/\sqrt{q\lambda_{E,\max}(n)}$. By the fact that $\forall \epsilon > 0$,

$$266 \quad (3.8) \quad P\left(h_n \|\beta_n - \theta\| > \tilde{M}\right) \leq P(\|\beta_n - \theta\| \geq \epsilon/2) + P\left(\tilde{M}/h_n < \|\beta_n - \theta\| < \epsilon/2\right),$$

267 we divide the proof into two steps: one is to prove $P(\|\beta_n - \theta\| \geq \epsilon/2) \leq \frac{\epsilon_1}{3}$ and the
 268 other is to prove $P\left(\tilde{M}/h_n < \|\beta_n - \theta\| < \epsilon/2\right) \leq \frac{2\epsilon_1}{3}$. Specifics are as follows.

269 *Step 1:* By Theorem 3.7, the probability $P(\|\beta_n - \theta\| \geq \epsilon/2)$ converges to zero,
 270 which means for any given $\epsilon_1 > 0$, there is a finite $N_1 \in \mathbb{N}_+$ such that for all $n > N_1$,

$$271 \quad (3.9) \quad P(\|\beta_n - \theta\| \geq \epsilon/2) \leq \epsilon_1/3.$$

272 *Step 2:* This step is to prove $P\left(\tilde{M}/h_n < \|\beta_n - \theta\| < \epsilon/2\right) \leq \frac{2\epsilon_1}{3}$. For each $n \in$
 273 \mathbb{N}_+ , denote $S_{j,n} = \{\beta : 2^{j-1} < h_n \|\beta - \theta\| < 2^j\}$ for $j \in \mathbb{Z}$. By Assumption (A2)(b),
 274 there is a finite $M_1 > 0$ and a finite $N_2 \in \mathbb{N}_+$ such that for all $n > N_2$,

$$275 \quad (3.10) \quad P(\lambda_{\min}(n) \leq M_1 d_n) = P(d_n \lambda_{\min}(n)^{-1} \geq M_1^{-1}) \leq \frac{\epsilon_1}{3}.$$

276 Denote $A_n = \{\omega : \lambda_{\min}(n) \leq M_1 d_n\}$. By the definition of $S_{j,n}$ and (3.10), we have

$$\begin{aligned} 277 & P(2^M/h_n < \|\beta_n - \theta\| < \epsilon/2) \\ 278 & \leq P(\{\omega : \beta_n \in S_{j,n}, \forall j \geq M+1, 2^{j+1} \leq \epsilon h_n\} \cap A_n^c) + P(A_n) \\ 279 \quad (3.11) & \leq \sum_{j \geq M+1, 2^{j+1} \leq \epsilon h_n} P(\{\omega : \beta_n \in S_{j,n}\} \cap A_n^c) + \frac{\epsilon_1}{3}. \end{aligned}$$

280 Since β_n is the minimum of $J_n(\beta)$, for any set A containing the point β_n , we have
 281 $\inf_{\beta \in A} (J_n(\beta) - J_n(\theta)) \leq 0$, which implies

$$282 \quad (3.12) \quad \{\omega : \beta_n \in A\} \subset \{\omega : \inf_{\beta \in A} (J_n(\beta) - J_n(\theta)) \leq 0\}.$$

283 Thus, by (3.12) and (3.11) we have

$$\begin{aligned} 284 & P(2^M/h_n < \|\beta_n - \theta\| < \epsilon/2) \\ 285 \quad (3.13) & \leq \frac{\epsilon_1}{3} + \sum_{j \geq M, 2^j \leq \epsilon h_n} P\left(\left\{\inf_{\beta \in S_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0\right\} \cap A_n^c\right). \end{aligned}$$

286 Next we consider the right hand of (3.13). Let $\beta = [\beta(1), \dots, \beta(p)]^T \in S_{j,n}$. Since
 287 $|\beta(j)|^\gamma > 0$ and $\theta(j) = 0$ for $j = q, q+1, \dots, p$, similar to (3.5), we have

$$\begin{aligned} 288 & J_n(\beta) - J_n(\theta) = (\beta - \theta)^T \sum_{i=1}^n (\varphi_i \varphi_i^T) (\beta - \theta) \\ 289 \quad (3.14) & + 2 \sum_{i=1}^n \varphi_i^T (\theta - \beta) w_{i+1} + \lambda_n \sum_{j=1}^q [|\beta(j)|^\gamma - |\theta(j)|^\gamma]. \end{aligned}$$

290 For the first term on the right hand of (3.14), by noting $\beta \in S_{j,n}$ and (3.10), for any
 291 $\omega \in A_n^c$, it follows that

$$\begin{aligned} 292 & (\beta - \theta)^T \sum_{i=1}^n (\varphi_i \varphi_i^T) (\beta - \theta) \geq \lambda_{\min}(n) \|\beta - \theta\|^2 \\ 293 \quad (3.15) & \geq \lambda_{\min}(n) 2^{2j-2} h_n^{-2} \geq M_1 d_n 2^{2j-2} h_n^{-2}. \end{aligned}$$

For the third term on the right hand of (3.14), by the mean value theorem, there exists ξ_j between $\beta(j)$ and $\theta(j)$ such that

$$\lambda_n \sum_{j=1}^q \left| |\beta(j)|^\gamma - |\theta(j)|^\gamma \right| = \lambda_n \gamma \sum_{j=1}^q |\xi_j|^{\gamma-1} |\beta(j) - \theta(j)|.$$

294 Since $\|\beta - \theta\| < \epsilon/2$, there is some constant $C_1 > 0$ such that $|\xi_j|^{\gamma-1} < C_1$. Thus,

$$\begin{aligned} 295 \quad & \lambda_n \sum_{j=1}^q \left| |\beta(j)|^\gamma - |\theta(j)|^\gamma \right| \leq C_1 \lambda_n \gamma \sum_{j=1}^q |\beta(j) - \theta(j)| \\ 296 \quad (3.16) \quad & \leq C_1 \lambda_n \gamma \sqrt{q} \|\beta - \theta\| \leq C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1}. \end{aligned}$$

297 Then, it follows from (3.16) that

$$298 \quad (3.17) \quad \lambda_n \sum_{j=1}^q \left[|\beta(j)|^\gamma - |\theta(j)|^\gamma \right] \geq -C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1}.$$

299 Hence, by (3.14), (3.15) and (3.17), we have for any $\omega \in A_n^c$,

$$300 \quad (3.18) \quad J_n(\beta) - J_n(\theta) \geq M_1 d_n 2^{2j-2} h_n^{-2} - C_1 \lambda_n \gamma 2^j h_n^{-1} - \sup_{\beta \in S_{j,n}} 2 \left| \sum_{i=1}^n \varphi_i^T (\theta - \beta) w_{i+1} \right|.$$

301 When $\inf_{\beta \in S_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0$, by (3.18), for any $\omega \in A_n^c$, the following inequality
302 holds

$$303 \quad (3.19) \quad \sup_{\beta \in S_{j,n}} 2 \left| \sum_{i=1}^n \varphi_i^T (\theta - \beta) w_{i+1} \right| \geq M_1 d_n 2^{2j-2} h_n^{-2} - C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1}.$$

304 By Assumption (A3)(b), we have $\frac{\lambda_n \sqrt{q} 2^j h_n^{-1}}{d_n 2^{2j-2} h_n^{-2}} = \frac{\lambda_n}{2^{j-2} \sqrt{\lambda_{E, \max}(n)}} \xrightarrow{n \rightarrow \infty} 0$. Then, it
305 follows that $M_1 d_n 2^{2j-2} h_n^{-2} > C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1}$ for all $n > N_3$ with N_3 being some
306 positive integer. Therefore, by (3.19) and Markov inequality, we have

$$\begin{aligned} 307 \quad & \mathbb{P} \left(\left\{ \inf_{\beta \in S_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0 \right\} \cap A_n^c \right) \\ 308 \quad & \leq \mathbb{P} \left(\sup_{\beta \in S_{j,n}} 2 \left| \sum_{i=1}^n \varphi_i^T (\theta - \beta) w_{i+1} \right| \geq M_1 d_n 2^{2j-2} h_n^{-2} - C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1} \right) \\ 309 \quad (3.20) \quad & \leq \frac{E \sup_{\beta \in S_{j,n}} 2 \left| \sum_{i=1}^n \varphi_i^T (\theta - \beta) w_{i+1} \right|}{M_1 d_n 2^{2j-2} h_n^{-2} - C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1}}. \end{aligned}$$

310 In addition, by Assumption (A1), we further assume that $E(w_{k+1}^2 | \mathcal{F}_k) = \sigma_k^2 \leq \bar{\sigma}^2$
311 with $\bar{\sigma}$ being some constant. Then, by the definition of $S_{j,n}$, Jensen's inequality, and
312 Cauchy-Schwarz inequality, we have

$$\begin{aligned} 313 \quad & E \sup_{\beta \in S_{j,n}} 2 \left| \sum_{i=1}^n \varphi_i^T (\theta - \beta) w_{i+1} \right| \leq 2 \sqrt{E \sup_{\beta \in S_{j,n}} \|\beta - \theta\|^2 \left\| \sum_{i=1}^n \varphi_i^T w_{i+1} \right\|^2} \\ 314 \quad (3.21) \quad & \leq 2^{j+1} h_n^{-1} \sqrt{E \left[\sum_{i=1}^n \varphi_i^T w_{i+1} \sum_{i=1}^n \varphi_i w_{i+1} \right]}. \end{aligned}$$

315 Noting Assumption (A1), we have

$$\begin{aligned}
316 \quad E \left[\sum_{i=1}^n \varphi_i^T w_{i+1} \sum_{i=1}^n \varphi_i w_{i+1} \right] &= E \left[\sum_{i=1}^n \varphi_i^T \varphi_i w_{i+1}^2 \right] = E \left[\sum_{i=1}^n E([\varphi_i^T \varphi_i w_{i+1}^2] \mid \mathcal{F}_i) \right] \\
317 \quad (3.22) \quad &\leq \bar{\sigma}^2 E \sum_{i=1}^n \varphi_i^T \varphi_i \leq \bar{\sigma}^2 \text{tr} \left(E \sum_{i=1}^n \varphi_i \varphi_i^T \right) \leq \bar{\sigma}^2 p \lambda_{E, \max}(n).
\end{aligned}$$

318 Therefore, from (3.20) and (3.22) it follows that

$$\begin{aligned}
319 \quad &\mathbb{P} \left(\left\{ \inf_{\beta \in \mathcal{S}_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0 \right\} \cap A_n^c \right) \\
320 \quad &\leq \frac{2^{j+1} \bar{\sigma} \sqrt{p} h_n^{-1} \lambda_{E, \max}(n)^{1/2}}{M_1 d_n 2^{2j-2} h_n^{-2} - C_1 \lambda_n \gamma \sqrt{q} 2^j h_n^{-1}} \leq \frac{2 \bar{\sigma}}{M_1 2^{j-2} - \frac{C_1 \lambda_n \gamma}{\lambda_{E, \max}(n)^{1/2}}}.
\end{aligned}$$

321 By Assumption (A3)(b), there is a finite $N_4 \in \mathbb{N}_+$ such that for all $n > N_4$,

$$322 \quad \mathbb{P} \left(\left\{ \inf_{\beta \in \mathcal{S}_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0 \right\} \cap A_n^c \right) \leq \frac{\bar{\sigma}}{M_1 2^{j-4}},$$

323 which leads to

$$324 \quad (3.23) \quad \sum_{j \geq M, 2^j \leq \epsilon h_n} \mathbb{P} \left(\left\{ \inf_{\beta \in \mathcal{S}_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0 \right\} \cap A_n^c \right) \leq \sum_{j \geq M} \frac{\bar{\sigma}}{M_1 2^{j-4}} \leq \frac{\bar{\sigma}}{M_1} 2^{-(M-5)}.$$

325 Therefore, for the given ϵ_1 , there is a finite M_2 and N_5 such that for all $n > N_5$,

$$326 \quad (3.24) \quad \sum_{j \geq M_2, 2^j \leq \epsilon h_n} \mathbb{P} \left(\left\{ \inf_{\beta \in \mathcal{S}_{j,n}} (J_n(\beta) - J_n(\theta)) \leq 0 \right\} \cap A_n^c \right) \leq \frac{\epsilon_1}{3}.$$

327 Thus, from (3.8), (3.9), (3.13) and (3.24), letting $\tilde{M} = 2^{M_2}$ and $N = \max\{N_1, \dots, N_5\}$,
328 we have for all $n > N$, $\mathbb{P} \left(h_n \|\beta_n - \theta\| > \tilde{M} \right) \leq \epsilon_1$. This completes the proof. \square

329 *Remark 3.9.* From Equation (3.7), it can be obtained that the smaller the number
330 of non-zero elements of the parameter vector θ , the faster the convergence rate
331 and thus the better the identification performance. This is further verified by the
332 simulation Example 1 in Section 6.

333 Based on Lemma 3.8, we give the following theorem demonstrating the set convergence
334 in probability. Different from Theorem 3.7, the following theorem further
335 states that the probability that the proposed algorithm correctly selects the non-zero
336 elements of the unknown sparse parameter vector converges to one.

337 **THEOREM 3.10.** *Let $\beta_n = (\beta_{1n}^T, \beta_{2n}^T)^T$ with $\beta_{1n} \in \mathbb{R}^q$ and $\beta_{2n} \in \mathbb{R}^{p-q}$ being
338 the vectors composed by the first q elements and the last $p-q$ elements of β_n , and
339 $\theta = (\theta_{10}^T, 0_{p-q}^T)^T$ with $\theta_{10} \in \mathbb{R}^q$. If Assumptions (A1)-(A3) hold, then we have the
340 set convergence of the estimates with probability tending to one, i.e., $\lim_{n \rightarrow \infty} P(\beta_{2n} =$
341 $0_{p-q}) = 1$.*

342 *Proof.* Let $t_n = \frac{\sqrt{\lambda_{E, \max}(n)}}{d_n}$. Denote the estimate $\beta_n = (\beta_{1n}^T, \beta_{2n}^T)^T$ as $\beta_{1n} =$
343 $\theta_{10} + t_n u_{1n}$, $\beta_{2n} = t_n u_{2n}$, where $u_{1n} \in \mathbb{R}^q$ and $u_{2n} \in \mathbb{R}^{p-q}$. In addition, denote

$$344 \quad (3.25) \quad \sum_{k=1}^n \varphi_k \varphi_k^T = \begin{bmatrix} \Phi_n^{(11)} & \Phi_n^{(12)} \\ \Phi_n^{(21)} & \Phi_n^{(22)} \end{bmatrix} \quad \text{and} \quad \varphi_k = \begin{bmatrix} \varphi_k^{(1)} \\ \varphi_k^{(2)} \end{bmatrix},$$

345 where $\Phi_n^{(11)} \in \mathbb{R}^{q \times q}$, $\varphi_k^{(q)} \in \mathbb{R}^d$ and others are with compatible dimensions. For any
 346 given $C > 0$, define

$$347 \quad (3.26) \quad B_n = \{\omega : \beta_n \in \{\beta : \|\beta - \theta\| \leq t_n C\}\}, \quad D_n = \{\omega : \beta_{2n} = 0\}.$$

348 Then, by (3.26) we have $\|u_{1n}\| \leq C$ and $\|u_{2n}\| \leq C$ for all $\omega \in B_n$. Next, we prove
 349 that for any given $\epsilon > 0$, there is $N \in \mathbb{N}_+$ such that $\mathbb{P}(D_n) \geq 1 - \epsilon$ for all $n > N$.
 350 In the following, we consider the estimate sequence $\{\beta_n\}_{n \geq 1}$ on a fixed sample path
 351 $\omega \in B_n$. Direct calculation for (3.2) leads to

$$352 \quad J_n(\beta_n) = \sum_{i=1}^n w_{i+1}^2 + (\beta_n - \theta)^T \sum_{i=1}^n (\varphi_i \varphi_i^T) (\beta_n - \theta)$$

$$353 \quad + 2 \sum_{i=1}^n \varphi_i^T (\theta - \beta_n) w_{i+1} + \lambda_n \sum_{j=1}^p |\beta_n(j)|^\gamma.$$

354 Then, we can obtain

$$355 \quad J_n(\theta_{10} + t_n u_{1n}, t_n u_{2n}) - J_n(\theta_{10} + t_n u_{1n}, 0)$$

$$356 \quad = t_n^2 \sum_{i=1}^n \left(\varphi_i^{(2)T} u_{2n} \right)^2 + 2t_n^2 \sum_{i=1}^n \left(\varphi_i^{(1)T} u_{1n} \right) \left(\varphi_i^{(2)T} u_{2n} \right)$$

$$357 \quad (3.27) \quad - 2t_n \sum_{i=1}^n w_{i+1} \left(\varphi_i^{(2)T} u_{2n} \right) + \lambda_n t_n^\gamma \sum_{j=1}^{p-q} |u_{2n}(j)|^\gamma.$$

358 For the first two terms on the right hand of (3.27), we have

$$359 \quad (3.28) \quad t_n^2 \sum_{i=1}^n \left(\varphi_i^{(2)T} u_{2n} \right)^2 + 2t_n^2 \sum_{i=1}^n \left(\varphi_i^{(1)T} u_{1n} \right) \left(\varphi_i^{(2)T} u_{2n} \right)$$

$$360 \quad \geq t_n^2 \sum_{i=1}^n \left(\varphi_i^{(2)T} u_{2n} \right)^2 - t_n^2 \sum_{i=1}^n \left[\left(\varphi_i^{(1)T} u_{1n} \right)^2 + \left(\varphi_i^{(2)T} u_{2n} \right)^2 \right] = -t_n^2 \sum_{i=1}^n \left(\varphi_i^{(1)T} u_{1n} \right)^2.$$

361 By Markov inequality and noting that $\lambda_{\max}\{E\Phi_n^{(11)}\} \leq \lambda_{E,\max}(n)$, for the above given
 362 ϵ , letting $M_1 = \frac{3}{\epsilon}$, we have

$$363 \quad (3.29) \quad \mathbb{P} \left(t_n^2 \sum_{i=1}^n \left(\varphi_i^{(1)T} u_{1n} \right)^2 \geq M_1 \lambda_{E,\max}(n) t_n^2 C^2 \right) \leq \frac{\epsilon E \left(t_n^2 \sum_{i=1}^n \left(\varphi_i^{(1)T} u_{1n} \right)^2 \right)}{3 \lambda_{E,\max}(n) t_n^2 C^2} \leq \epsilon/3.$$

364 Hence, it follows $\mathbb{P}(E_n^c) \leq \epsilon/3$, where E_n is denoted as

$$365 \quad (3.30) \quad E_n = \left\{ \omega : t_n^2 \sum_{i=1}^n \left(\varphi_i^{(1)T} u_{1n} \right)^2 \leq M_1 t_n^2 C^2 \lambda_{E,\max}(n) \right\}.$$

366 For the third term on the right hand of (3.27), similar to (3.21) and (3.22), noting
 367 that $\lambda_{\max}\{E\Phi_n^{(22)}\} \leq \lambda_{E,\max}(n)$ and $\|u_{2n}\| \leq C$, we can get

$$368 \quad E \left| \sum_{i=1}^n w_{i+1} \left(\varphi_i^{(2)T} u_{2n} \right) \right| \leq \left(E \left| \sum_{i=1}^n w_{i+1} \left(\varphi_i^{(2)T} u_{2n} \right) \right|^2 \right)^{1/2}$$

$$369 \quad (3.31) \quad \leq C \bar{\sigma} \lambda_{\max}\{E\Phi_n^{(22)}\}^{1/2} \leq C \bar{\sigma} \sqrt{\lambda_{E,\max}(n)},$$

370 where $\mathbb{E}(w_{k+1}^2 | \mathcal{F}_k) \leq \bar{\sigma}^2$ with $\bar{\sigma}$ being some constant by Assumption (A1). By Markov
 371 inequality, for the above given ϵ , letting $M_2 = \frac{3C\bar{\sigma}}{\epsilon}$, it follows from (3.31) that

$$372 \quad (3.32) \quad \mathbb{P} \left(\left| \sum_{i=1}^n w_{i+1} \left(\varphi_i^{(2)T} u_{2n} \right) \right| \geq M_2 \sqrt{\lambda_{E, \max}(n)} \right) \leq \frac{E \left| \sum_{i=1}^n w_{i+1} \left(\varphi_i^{(2)T} u_{2n} \right) \right|}{M_2 \sqrt{\lambda_{E, \max}(n)}} \leq \epsilon/3.$$

373 Denote

$$374 \quad (3.33) \quad F_n = \left\{ \omega : - \sum_{i=1}^n w_{i+1} \left(\varphi_i^{(2)T} u_{2n} \right) \geq -M_2 \lambda_{E, \max}^{1/2}(n) \right\}.$$

375 Then, from (3.32) it follows $\mathbb{P}(F_n^c) \leq \epsilon/3$. For the last term on the right hand of
 376 (3.27), noting that $\left[\sum_{j=1}^{p-q} |u_{2n}(j)|^\gamma \right]^{2/\gamma} \geq \sum_{j=1}^{p-q} |u_{2n}(j)|^2 = \|u_{2n}\|^2$, we have

$$377 \quad (3.34) \quad \lambda_n t_n^\gamma \sum_{j=1}^{p-q} |u_{2n}(j)|^\gamma \geq \|u_{2n}\|^\gamma \lambda_n t_n^\gamma.$$

378 For all $\omega \in E_n \cap F_n$, from (3.28), (3.30), (3.33) and (3.34), we have

$$379 \quad J_n(\theta_{10} + t_n u_{1n}, t_n u_{2n}) - J_n(\theta_{10} + t_n u_{1n}, 0) \geq \\
 380 \quad (3.35) \quad -M_1 t_n^2 C^2 \lambda_{E, \max}(n) + \|u_{2n}\|^\gamma \lambda_n t_n^\gamma - 2t_n M_2 \lambda_{E, \max}^{1/2}(n).$$

381 By Assumption (A3)(c), and noting that $\lambda_{E, \max}(n)/d_n \rightarrow 0$, we have

$$382 \quad \frac{\lambda_n t_n^\gamma}{t_n^2 \lambda_{E, \max}(n)} = \frac{\lambda_n d_n^{2-\gamma}}{\lambda_{E, \max}(n)^{2-\frac{1}{2}\gamma}} \xrightarrow{n \rightarrow \infty} \infty, \\
 383 \quad \frac{\lambda_n t_n^\gamma}{t_n \sqrt{\lambda_{E, \max}(n)}} = \frac{\lambda_n d_n^{2-\gamma}}{\lambda_{E, \max}(n)^{2-\frac{1}{2}\gamma}} \frac{\lambda_{E, \max}(n)}{d_n} \xrightarrow{n \rightarrow \infty} \infty.$$

384 Therefore, from (3.35), if $\|u_{2n}\| > 0$, then there is a finite $\tilde{N} \in \mathbb{N}_+$ such that $J_n(\beta_n) -$
 385 $J_n(\theta_{10} + t_n u_{1n}, 0) > 0$, $\forall n > \tilde{N}$, which contradicts $\beta_n = \underset{\beta}{\operatorname{argmin}} J_n(\beta)$. Thus, for any
 386 $\omega \in E_n \cap F_n$, there is a finite N_1 such that $\beta_{2n} = t_n u_{2n} = 0$, $\forall n > N_1$. Therefore,
 387 from (3.26) it follows $B_n \cap E_n \cap F_n \subset D_n \cap E_n \cap F_n$, $\forall n > N_1$. In addition, by
 388 Lemma 3.8, for the above given ϵ , there is an $N_2 \in \mathbb{N}_+$ such that $P(B_n^c) \leq \epsilon/3$ for
 389 all $n > N_2$. Hence, combing the results above (3.30) and below (3.33), and letting
 390 $N = \max\{N_1, N_2\}$, we have that for all $n > N$,

$$391 \quad P(\beta_{2n} = 0) = P(D_n) \geq P(D_n \cap E_n \cap F_n) = 1 - P(D_n^c \cup E_n^c \cup F_n^c) \\
 392 \quad \geq 1 - P(B_n^c) - P(E_n^c) - P(F_n^c) \geq 1 - \epsilon.$$

393 This completes the proof. \square

394 Using the central limit theorem, we immediately give the asymptotic normality
 395 of the estimated non-zero parameters below.

396 **THEOREM 3.11.** *Assume for each n that there is a non-random positive definite*
 397 *symmetric matrix R_n such that*

$$398 \quad (3.36) \quad R_n^{-1} \Phi_n^{(11)} \xrightarrow{P} \mathbf{I}_p, \quad \max_{1 \leq k \leq n} \|R_n^{-1/2} \varphi_k^{(1)}\| \xrightarrow{P} 0, \quad \text{and}$$

$$399 \quad (3.37) \quad \lim_{k \rightarrow \infty} \mathbb{E}(w_{k+1}^2 | \mathcal{F}_k) = \sigma^2, \quad \text{a.s. for some constant } \sigma,$$

400 where $\varphi_k^{(1)}$ and $\Phi_n^{(11)}$ are defined in (3.25). Denote the estimate $\beta_n = (\beta_{1n}^T, \beta_{2n}^T)^T$
 401 and $\theta = [\theta_{10}^T, 0_{p-q}^T]^T$. For any non-random $\alpha_n \in \mathbb{R}^q$ satisfying $\|\alpha_n\| \leq 1$, let $s_n^2 =$
 402 $\sigma^2 \alpha_n^T R_n^{-1} \alpha_n$. If Assumptions (A1)-(A3) hold, then

$$403 \quad (3.38) \quad s_n^{-1} \alpha_n^T (\beta_{1n} - \theta_{10}) = s_n^{-1} \sum_{k=1}^n \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \varphi_k^{(1)} w_{k+1} + o_p(1) \xrightarrow{d} N(0, 1),$$

404 where \xrightarrow{d} denotes convergence in distribution and $N(0, 1)$ denotes the standard normal
 405 distribution.

406 *Remark 3.12.* The existence of a non-random matrix R_n satisfying conditions
 407 (3.36) in Theorem 3.11 can be regarded as an stability assumption of the matrix
 408 $\Phi_n^{(11)}$. Moreover, this assumption is necessary for asymptotic normality and one can
 409 refer to Example 3 in [21] in which asymptotic normality fails to hold in the absence of
 410 (3.36). Besides, R_n can be selected as $\Phi_n^{(11)}$ if $\{\varphi_k^{(1)}\}$ is determined sequence; R_n can
 411 be selected as $nE\varphi_n^{(1)}\varphi_n^{(1)T}$ if $\varphi_n^{(1)}\varphi_n^{(1)T}$ is a stationary and ergodic random sequence
 412 with positive covariance matrix [39].

413 *Proof.* Denote $J_n(\beta)$ in (3.2) as $J_n(\beta) = J_n(\beta_1, \beta_2)$ with $\beta_1 \in \mathbb{R}^q$. By Theorem
 414 3.7, we have $\|\beta_n - \theta\| \xrightarrow[n \rightarrow \infty]{} 0$ a.s. Since each component of θ_{10} is not equal to zero,
 415 when n is sufficiently large, each element of β_{1n} stays away from zero. Noting that
 416 the estimate $\beta_n = (\beta_{1n}^T, \beta_{2n}^T)^T$ is the minimum of $J_n(\beta)$, when n is sufficiently large,
 417 we have $\frac{\partial}{\partial \beta_1} J_n(\beta_{1n}, \beta_{2n}) = 0$, which implies

$$418 \quad (3.39) \quad -2 \sum_{k=1}^n \left(y_{k+1} - \beta_{1n}^T \varphi_k^{(1)} - \beta_{2n}^T \varphi_k^{(2)} \right) \varphi_k^{(1)} + \lambda_n \gamma \text{vec} \left(\text{sign}(\beta_{1n}(j)) |\beta_{1n}(j)|^{\gamma-1} \right) \Big|_{j=1}^q = 0.$$

419 From (2.1) and noting that $\theta = [\theta_{10}^T, 0_{1 \times (p-q)}]^T$, it follows $y_{k+1} - \theta_{10}^T \varphi_k^{(1)} = w_{k+1}$.
 420 Then, by (3.39) we get

$$421 \quad (3.40) \quad \sum_{k=1}^n \varphi_k^{(1)} \varphi_k^{(1)T} (\beta_{1n} - \theta_{10})$$

$$422 \quad = - \sum_{k=1}^n \beta_{2n}^T \varphi_k^{(2)} \varphi_k^{(1)} + \sum_{k=1}^n \varphi_k^{(1)} w_{k+1} - \frac{1}{2} \lambda_n \gamma \text{vec} \left(\text{sign}(\beta_{1n}(j)) |\beta_{1n}(j)|^{\gamma-1} \right) \Big|_{j=1}^q.$$

423 Thus, direct calculation from (3.40) leads to

$$424 \quad s_n^{-1} \alpha_n^T (\beta_{1n} - \theta_{10}) = -s_n^{-1} \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \sum_{k=1}^n \beta_{2n}^T \varphi_k^{(2)} \varphi_k^{(1)} + s_n^{-1} \sum_{k=1}^n \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \varphi_k^{(1)} w_{k+1}$$

$$425 \quad (3.41) \quad - \frac{1}{2} \lambda_n \gamma s_n^{-1} \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \text{vec} \left(\text{sign}(\beta_{1n}(j)) |\beta_{1n}(j)|^{\gamma-1} \right) \Big|_{j=1}^q.$$

426 For the first term on the right hand of (3.41), by Theorem 3.10 that $\lim_{n \rightarrow \infty} P(\beta_{2n} =$
 427 $0) = 1$, we have

$$428 \quad (3.42) \quad \lim_{n \rightarrow \infty} P \left(s_n^{-1} \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \sum_{k=1}^n \beta_{2n}^T \varphi_k^{(2)} \varphi_k^{(1)} = 0 \right) = 1.$$

429 For the last term on the right hand of (3.41), since $\beta_{1n} \rightarrow \theta_{10}$, there is a constant C
 430 such that $|\beta_{1n}(j)| \leq C$ for $j = 1, \dots, q$. By Assumption (A3)(a), we have

$$431 \quad \left| \lambda_n \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \text{vec} \left(\text{sign}(\beta_{1n}(j)) |\beta_{1n}(j)|^{\gamma-1} \right) \right|_{j=1}^q$$

$$432 \quad (3.43) \quad \leq \lambda_n \lambda_{\min}(n)^{-1} q^{1/2} C^{\gamma-1} = o_p(1),$$

433 which together with (3.41) and (3.42) gives

$$434 \quad (3.44) \quad s_n^{-1} \alpha_n^T (\beta_{1n} - \theta_{10}) = s_n^{-1} \sum_{k=1}^n \alpha_n^T \left(\Phi_n^{(11)} \right)^{-1} \varphi_k^{(1)} w_{k+1} + o_p(1).$$

435 In view of (3.36) and (3.44), to prove (3.38), we need only to show that

$$436 \quad (3.45) \quad s_n^{-1} \sum_{k=1}^n \alpha_n^T R_n^{-1} \varphi_k^{(1)} w_{k+1} \xrightarrow{d} N(0, 1).$$

437 Similar to [21], the desired conclusion (3.45) can be obtained by making use of a
 438 martingale central limit theorem of [7]. \square

439 **3.3. Comparison of Algorithm 3.1 with related methods.** In this part,
 440 we compare the sparse identification Algorithm 3.1 with Information Criterion-based
 441 variable selection [1, 32], LASSO [33], and bridge estimate [18].

442 *Comparison with variable selection and order estimation based on information*
 443 *criterion.* The variable selection problem aims to select a subset of relevant variables
 444 used in model construction. The usual approach is to select the optimal one from
 445 a set of reasonable models under some importance criteria, many of which contain
 446 measures of accuracy and the penalized term by the number of selected variables,
 447 for instance, AIC [1] and BIC [32] for stationary time series. Algorithm 3.1 in this
 448 paper not only fulfills the task of variable selection but also estimates the parameters
 449 corresponding to the selected variables. Moreover, compared with order estimation
 450 methods for stochastic systems such as control information criterion (CIC) [5, 15], the
 451 algorithm in this paper solves the problem as well, and furthermore, non-contributing
 452 variables within the order can also be selected out.

Comparison with LASSO and bridge estimate. Compared with the LASSO, Al-
 gorithm 3.1 does not require additional conditions; and compared with the bridge
 estimate, Algorithm 3.1 can be applied to general observations. In a typical setup,
 the sparsity problem can be described as follows [37]: Given a $n \times p$ matrix Ψ_n , and
 a procedure of generating an observation such as

$$Y = \Psi_n \theta + W$$

with $Y = [y_1, \dots, y_n]^T$, $\Psi_n = [\varphi_0, \dots, \varphi_{n-1}]^T$ and $W = [w_1, \dots, w_n]$, we are asked to
 recover θ from the observation Y such that θ is of the sparsest structure. The problem
 can be solved by the following regularization method:

$$\min_{\theta \in \mathbb{R}^p} \{ \|Y - \Psi_n \theta\|^2 + \lambda_n \|\theta\|_\nu^\nu \},$$

where $\nu > 0$ and $\|x\|_\nu$ is defined by $\|\theta\|_\nu = \sqrt[\nu]{\sum_{i=1}^p |\theta(i)|^\nu}$. The LASSO (for $\nu = 1$),
 the bridge estimate (for $\nu > 0$), and the Algorithm 3.1 (for $0 < \nu < 1$) in this paper all

fall into this category, but Ψ_n in Algorithm 3.1 can be stochastic, whereas the others are deterministic. We then compare the application scope of these three algorithms. For LASSO, denote

$$\Phi_n = \sum_{k=1}^n \varphi_k \varphi_k^T = \begin{bmatrix} \Phi_n^{11} & \Phi_n^{12} \\ \Phi_n^{21} & \Phi_n^{22} \end{bmatrix}$$

with $\Phi_n^{11} \in \mathbb{R}^{q \times q}$ and $\beta_n = (\beta_{1n}^T, \beta_{2n}^T)^T$. [40] gave sufficient conditions for the set convergence in probability of the LASSO estimate: (a) $\frac{1}{n} \Phi_n \rightarrow \Phi$ with Φ being a positive definite matrix; (b) the following strong irrepresentable condition holds:

$$\left| \Phi_n^{21} (\Phi_n^{11})^{-1} \text{sign}(\beta_{1n}) \right| \leq \mathbf{1}_{p-q} - \eta,$$

with $\mathbf{1}_{p-q}$ being a $(p-q) \times 1$ vector of 1's, $\eta > 0$ and the inequality holding element-wise; and (c) λ_n is chosen as $\lambda_n = n^\alpha$ with $\frac{1}{2} < \alpha < 1$. Algorithm 3.1 of this paper can also reach set convergence while covering condition (a) as a special case without requiring the strong irrepresentable condition (b).

For bridge estimate, the conditions for the consistency of the estimates given by [18] are: (a) $\frac{1}{n} \Phi_n \rightarrow \Phi$ with Φ being a positive definite matrix; (b) $\lambda_n n^{-1/2} \rightarrow 0$ and $\lambda_n^2 n^{-\gamma} \rightarrow 0$. This result is consistent with the result of Algorithm 3.1 when $C_1 n \leq \lambda_{\min}(n) \leq \lambda_{\max}(n) \leq C_2 n$ for some constants C_1 and C_2 . In addition, the theoretical results of Algorithm 3.1 go further and can be adapted to non-persistent excitation cases, in particular, the data series $\{\varphi_k, y_{k+1}\}_{k \geq 1}$ can be generated by feedback control where φ_k may be stochastic.

4. Weighted L_γ regularization algorithm and its properties. LASSO is a popular technique for simultaneous estimation and variable selection. However, in some cases, LASSO is inconsistent for variable selection. [28] showed the conflict between the optimal prediction and consistent variable selection in LASSO. To address this issue, [43] proposed a new version of the LASSO, the adaptive LASSO, in which adaptive weights were used to penalize different parameters in the L_1 penalty. [42] extended this result to general observation cases. Inspired by the improvement of the convergence properties of LASSO with this technique, in order to extend the scope of application and improve the performance of the L_γ penalty, in this section, we present a two-step algorithm with adaptively weighted L_γ ($0 < \gamma \leq 1$) penalty term. The algorithm is more broadly applicable and has better convergence properties.

4.1. Weighted L_γ regularization algorithm.

Assumption. Given constants γ and μ satisfying $0 < \gamma \leq 1$ and $\mu > 0$. To proceed, we first introduce the assumptions to be used for the theoretical analysis of the weighted L_γ regularization algorithm.

(B1) For the maximal and minimal eigenvalues of $\sum_{k=1}^n \varphi_k \varphi_k^T$ and the positive sequence $\{\lambda_n\}_{n \geq 1}$, it holds,

$$(a) \left(\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} \right)^{1 - \frac{\gamma}{2} + \frac{\mu}{2}} \frac{\lambda_{\max}(n)}{\lambda_n} \xrightarrow{n \rightarrow \infty} 0 \text{ a.s.}$$

$$(b) \frac{\log \lambda_{\max}(n)^{\frac{\mu}{2}}}{\lambda_{\min}(n)^{1 - \frac{\gamma}{2} + \frac{\mu}{2}}} \frac{\lambda_{\max}(n)}{\lambda_n^{\frac{\gamma}{2}}} \xrightarrow{n \rightarrow \infty} 0 \text{ a.s.} \quad (c) \frac{\log \lambda_{\max}(n)^{\frac{1}{2} + \frac{\mu}{2}}}{\lambda_{\min}(n)^{\frac{1}{2} - \frac{\gamma}{2} + \frac{\mu}{2}}} \frac{\lambda_{\max}(n)^{\frac{1}{2}}}{\lambda_n^{\frac{1}{2} + \frac{\mu}{2}}} \xrightarrow{n \rightarrow \infty} 0 \text{ a.s.}$$

The adaptive sparse identification algorithm is proposed in Algorithm 4.1.

Algorithm 4.1 Weighted L_γ regularization.

Step 0 (Initialization). For given $0 < \gamma \leq 1$ and $\mu > 0$, choose a positive sequence $\{\lambda_n\}_{n \geq 1}$ satisfying Assumption (B1).

Step 1 (LS Estimation). Based on $\{y_{k+1}, \varphi_k\}_{k=1}^n$, compute the estimator:

$$\theta_{n+1} = \left(\sum_{k=1}^n \varphi_k \varphi_k^T \right)^{-1} \left(\sum_{k=1}^n \varphi_k y_{k+1} \right).$$

Let $\theta_{n+1} = [\theta_{n+1}(1), \dots, \theta_{n+1}(p)]^T$, and for $1 \leq j \leq p$, define

$$\hat{\theta}_{n+1}(j) = \theta_{n+1}(j) + \text{sign}(\theta_{n+1}(j)) \sqrt{\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)}}.$$

Step 2 (Sparse Optimization with L_γ penalty). With λ_n and $\hat{\theta}_{n+1}$, optimize the objective function $\hat{J}_n(\beta) = \sum_{k=1}^n (y_{k+1} - \beta^T \varphi_k)^2 + \lambda_n \sum_{j=1}^p \frac{1}{|\hat{\theta}_{n+1}(j)|^\mu} |\beta(j)|^\gamma$ and obtain

$$(4.1) \quad \hat{\beta}_n = \left[\hat{\beta}_n(1), \dots, \hat{\beta}_n(p) \right]^T = \underset{\beta}{\text{argmin}} \hat{J}_n(\beta)$$

$$(4.2) \quad \hat{A}_n^* = \left\{ j = 1, \dots, p \mid \hat{\beta}_n(j) = 0 \right\}.$$

480 *Remark 4.1.* We discuss the choice of λ_n in the Algorithm 4.1. If we as-
 481 sume $\frac{\lambda_{\max}(n)}{\lambda_{\min}(n)} \left(\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} \right)^{\mu/2} \rightarrow 0$, a.s., then Assumption (B1) can be simplified to
 482 $\lambda_n = o(\lambda_{\min}(n))$ and $\lambda_{\max}(n) \left(\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} \right)^{\frac{\mu}{2}} = o(\lambda_n)$. Denote $a_n = \lambda_{\max}(n) \left(\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} \right)^{\frac{\mu}{2}}$
 483 and $b_n = \lambda_{\min}(n)$. Then, λ_n can be chosen as $\lambda_n = a_n^\eta b_n^{1-\eta}$ for any fixed $\eta \in (0, 1)$ sat-
 484 isfying Assumption (B1). Specifically, by noting that $\frac{a_n}{b_n} = \frac{\lambda_{\max}(n)}{\lambda_{\min}(n)} \left(\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} \right)^{\mu/2} \rightarrow$
 485 0 a.s., it follows that $\frac{\lambda_n}{b_n} = \left(\frac{a_n}{b_n} \right)^\eta \rightarrow 0$, and $\frac{a_n}{\lambda_n} = \left(\frac{a_n}{b_n} \right)^{1-\eta} \rightarrow 0$ a.s.

486 **4.2. Theoretical properties.** Recall that the parameter vector is assumed $\theta =$
 487 $[\theta(1), \dots, \theta(q), \theta(q+1), \dots, \theta(p)]^T$ with $\theta(i) \neq 0$ for $i = 1, \dots, q$, and $\theta(j) = 0$ for
 488 $j = q+1, \dots, p$. For the estimate $\hat{\beta}_n$ and \hat{A}_n^* generated by Algorithm 4.1, the almost
 489 sure convergence of $\hat{\beta}_n$ and the almost sure set convergence of \hat{A}_n^* are given in the
 490 following theorems.

THEOREM 4.2. *If Assumptions (A1), (A2)(a) and (A3)(a) hold, then*

$$\lim_{n \rightarrow \infty} \hat{\beta}_n(j) = \theta(j), \quad j = 1, \dots, q, \quad \text{a.s.}$$

491 *Proof.* The proof is similar to that of Theorem 3.7, and so, omitted here.

492 **THEOREM 4.3.** *If Assumptions (A1), (A3)(a) and (B1) hold, then there is an*
 493 *ω -space Ω_0 satisfying $P(\Omega_0) = 1$ and for any $\omega \in \Omega_0$, there is an integer $N_0(\omega)$ such*
 494 *that $\hat{A}_n^* = A^*$ for all $n \geq N_0(\omega)$.*

495 *Proof.* Combining the proof of Theorem 3.10 with the proof of Lemma 4 in [42]
 496 yields the theorem. \square

497 **4.3. Comparison of Algorithms 4.1 with related methods.** Noting Re-
 498 mark 3.4, Algorithm 4.1 is more likely to produce sparse solutions than the algorithm
 499 in [42] and adaptive LASSO [43]. Moreover, Algorithm 4.1 covers the results of the
 500 adaptive LASSO and the algorithm in [42]. Specifically, when $\mu = \gamma = 1$, by Remark
 501 4.1, Assumption (B1) is degenerated to $\frac{\lambda_{\max}(n)}{\lambda_{\min}(n)} \sqrt{\frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)}} \rightarrow 0$ a.s., which is consis-
 502 tent with Assumption (A3) in [42]. If one further assumes that $\lambda_{\max}(n) = O(n)$ and
 503 $\lambda_{\min}(n) = O(n)$ a.s., then the result is consistent with the adaptive LASSO.

504 5. Application to typical sciences.

505 **5.1. Structure selection for a class of NARX models.** This section ap-
 506 plies Algorithm 3.1 to the structure selection of the NARX models with finite basis
 507 functions. One class of NARX models [41] is the kernel regression model:

$$508 \quad (5.1) \quad y_{k+1} = \theta_N^T \varphi_{N,k} + w_{k+1},$$

509 where y_{k+1} is the output, $\varphi_{N,k} = [\varphi_1(x(k)), \dots, \varphi_m(x(k))]^T$ is the non-linear basis
 510 functions, $x(k)$ contains all past and current variables, $\theta_N = [c_1, \dots, c_m]^T$ is the cor-
 511 responding coefficient, $w_{k+1} \in \mathbb{R}$ is the noise and m is the number of basis functions.
 512 The objective of the NARX model structure selection is to select the contributing
 513 components from a large number of non-linear basis functions. Algorithm 3.1 can
 514 be applied directly to the model (5.1), and is more efficient in reducing the model
 515 size. Now we consider a special class of the NARX model: Hammerstein system as an
 516 example and give the corresponding theoretical results. The Hammerstein model con-
 517 sists of a static single-valued nonlinear element followed by a linear dynamic element,
 518 and can be described as:

$$519 \quad (5.2) \quad y_{k+1} = a_1 y_k + \dots + a_{n_y} y_{k+1-n_y} + b_1 f(u_k) + \dots + b_{n_u} f(u_{k+1-n_u}) + w_{k+1},$$

$$520 \quad f(u_k) = \sum_{j=1}^s d_j g_j(u_k),$$

521 where $\{g_j(\cdot)\}_{j=1}^s$ are the basis functions. The system (5.2) can be rewritten as the
 522 form (5.1) by denoting

$$523 \quad \theta_N = [a_1, \dots, a_{n_y}, (b_1 d_1), \dots, (b_1 d_s), \dots, (b_{n_u} d_1), \dots, (b_{n_u} d_s)]^T,$$

$$\varphi_{N,k} = [y_k, \dots, y_{k+1-n_y}, g_1(u_k), \dots, g_s(u_k), \dots, g_1(u_{k+1-n_u}), \dots, g_s(u_{k+1-n_u})]^T.$$

524 **Problem.** The structure selection problem of the Hammerstein system (5.2)
 525 is to select the contributing basis functions from the candidate full basis functions
 526 $\{g_j(\cdot)\}_{j=1}^s$ using the observed data $\{y_{k+1}, \varphi_{N,k}\}_{k=1}^n$.

527 Before presenting the results, we first give the following assumptions and the
 528 corresponding proposition.

- 529 **(C1)** $A(z) = 1 - a_1 z - \dots - a_{n_y} z^{n_y}$ is stable and $b_1^2 + \dots + b_{n_u}^2 \neq 0$;
 530 **(C2)** There is an interval $[a, b]$ such that $\{1, g_1(x), \dots, g_s(x)\}$ is linearly independent;
 531 **(C3)** The sequence $\{u_k\}_{k \geq 1}$ is i.i.d, independent of the noise $\{w_k\}_{k \geq 1}$, whose density
 532 function is positive and continuous on $[a, b]$ and $0 < \mathbb{E}g_i^2(u_k) < \infty$ for $1 \leq i \leq s$.

533 PROPOSITION 5.1. [41] *If the Hammerstein system (5.2) satisfies Assumptions*
 534 *(A1) and (C1)-(C3), then with $0 < c_1 < c_2$, $0 < c_3 < c_4$, we have*

$$535 \quad (5.3) \quad c_1 n \leq \lambda_{\max} \left\{ \sum_{k=1}^n \varphi_{N,k} \varphi_{N,k}^T \right\} \leq c_2 n, \quad c_3 n \leq \lambda_{\min} \left\{ \sum_{k=1}^n \varphi_{N,k} \varphi_{N,k}^T \right\} \leq c_4 n,$$

536

537 By use of Algorithm 3.1, we give the sparse estimate $\beta_{N,n}$ for the parameters in the
 538 system (5.2). Denote $\beta_{N,n} = [a_{1,n}, \dots, a_{n_y,n}, (b_1 d_1)_n, \dots, (b_1 d_s)_n, \dots, (b_{n_u} d_1)_n, \dots, (b_{n_u} d_s)_n]^T$
 539 and define $\chi = [\chi(1), \dots, \chi(s)]^T$ with $\chi(l) = \sum_{i=1}^{n_u} (b_i d_l)^2$. Then, the estimate of χ can
 540 be obtained by $\chi_n = [\chi_n(1), \dots, \chi_n(s)]^T$ with $\chi_n(l) = \sum_{i=1}^{n_u} (b_i d_l)_n^2$. Moreover, denote
 541 $D^* = \{l : d_l = 0, \text{ for } l = 1, \dots, s\}$, $D_n^* = \{l : \chi_n(l) = 0, \text{ for } l = 1, \dots, s\}$. Assumption
 542 (C1) guarantees that $\{l : \chi(l) = 0\} = D^*$, which implies D_n^* can be regraded as an
 543 estimate of D^* . Then, we give the theoretical results for the structure selection of the
 544 contributing basis functions in $\{g_j(\cdot)\}_{j=1}^s$.

545 THEOREM 5.2. *Take $\lambda_n = n^\alpha$ with $\frac{1}{2}\gamma < \alpha < \frac{1}{2}$. If Assumptions (A1) and (C1)-*
 546 *(C3) hold for the Hammerstein system (5.2), then $\lim_{n \rightarrow \infty} P(D_n^* = D^*) = 1$.*

547 *Proof.* From (5.3) in Proposition 5.1, we have that $\lambda_{\max} \left\{ \sum_{k=1}^n \varphi_{N,k} \varphi_{N,k}^T \right\} =$
 548 $O(n)$, $\lambda_{\min} \left\{ \sum_{k=1}^n \varphi_{N,k} \varphi_{N,k}^T \right\} = O(n)$ and $E \lambda_{\max} \left\{ \sum_{k=1}^n \varphi_{N,k} \varphi_{N,k}^T \right\} = O(n)$. More-
 549 over, we can choose $d_n = c_5 n$ with $c_5 > 0$. Thus, noticing $\lambda_n = n^\alpha$ with $\frac{1}{2}\gamma < \alpha < \frac{1}{2}$,
 550 we can verify that (A2)-(A3) hold for the regression model (5.1)-(5.2). Thus, by
 551 Theorem 3.10, the results follow directly. \square

552 **5.2. Sparse identification of linear feedback control systems.** This sec-
 553 tion applies Algorithm 3.1 to the sparse identification of the closed-loop systems using
 554 the self-tuning regulator (STR) control. Recall that the regressor is generally non-
 555 stationary and non-independent for linear feedback control systems [17]. The classical
 556 STR control, first proposed in [2], consists of an LS estimation algorithm for a linear
 557 stochastic dynamic system coupled online with a ‘‘least variance’’ control law. The
 558 goal of STR is to minimize the tracking error of the system with unknown parameters.
 559 Consider the following sparse ARX system:

$$560 \quad (5.4) \quad y_{k+1} = a_1 y_k + \dots + a_{n_y} y_{k+1-n_y} + b_1 u_k + \dots + b_{n_u} u_{k+1-n_u} + w_{k+1}.$$

561 where $y_{k+1} \in \mathbb{R}$ is the system output, $w_{k+1} \in \mathbb{R}$ is the system noise, $u_k \in \mathbb{R}$ is the
 562 feedback control, and a_1, \dots, a_{n_y} and b_1, \dots, b_{n_u} are the unknown sparse parameters.
 563 Denote

$$564 \quad A(z) = 1 - a_1 z - \dots - a_{n_y} z^{n_y}, \quad B(z) = b_1 + b_2 z + \dots + b_{n_u} z^{n_u-1},$$

$$\theta = [a_1, \dots, a_{n_y}, b_1, \dots, b_{n_u}]^T, \quad \varphi_k = [y_k, \dots, y_{k+1-n_y}, u_k, \dots, u_{k+1-n_u}]^T.$$

565 Let $\{y_k^*\}$ be the deterministic bounded reference signal or regulation signal. For the
 566 system (5.4), two problems need to be solved: first, to use the STR control to make
 567 the closed-loop system track the reference signal $\{y_k^*\}$; second, to select the zero
 568 parameters accurately and estimate the non-zero parameters asymptotically under
 569 the STR control.

570 For the control step, let the LS parameter estimate for the system be $\theta_{L,n} =$
 571 $[a_{1,n}, \dots, a_{n_y,n}, b_{1,n}, \dots, b_{n_u,n}]^T$. The Certainty Equivalence Principle [2] suggests an
 572 adaptive control defined as

$$573 \quad (5.5) \quad u_k^0 = \frac{1}{b_{1,k}} \{y_{k+1}^* + (b_{1,k} u_k - \theta_{L,k}^T \varphi_k)\}.$$

574 For the identification step, it is generally necessary to impose excitation conditions
 575 on the system, and thus, the control design (5.5) needs to be modified. Specifically,
 576 in order not to make the system worse after applying the excitation, the diminishing
 577 excitation technique is introduced and a zero-trending perturbation [14] is added to
 578 the control (5.5), i.e.,

$$579 \quad (5.6) \quad u_k = u_k^0 + \frac{\nu_k}{r_{k-1}^{\bar{\varepsilon}/2}}, \quad k \geq 1,$$

580 where $\{\nu_k\}$ is an i.i.d and bounded stochastic sequence satisfying $E(\nu_k) = 0$, $E(\nu_k^2) =$
 581 1 , $r_{k-1} = 1 + \sum_{i=1}^{k-1} \|\varphi_i\|^2$, $\bar{\varepsilon} \in \left(0, \frac{1}{2(\bar{n}_{yu}+1)}\right)$ and $\bar{n}_{yu} = \max\{n_y, n_u\} + n_y - 1$. Next,
 582 we give the assumptions for (5.4) and the stability and optimality in Proposition 5.3.

583 **(D1)** The noise $\{w_k\}$ satisfies $\lim_{k \rightarrow \infty} \frac{1}{k} \sum_{j=1}^k w_j^2 = R > 0$ a.s.;

584 **(D2)** The system is of minimum phase, i.e., $B(z) \neq 0, \forall |z| \leq 1$;

585 **(D3)** $|a_{n_y}| + |b_{n_u}| \neq 0$.

PROPOSITION 5.3. [14] *If Assumptions (A1) and (D1)-(D3) hold, then the model (5.4) with the attenuating excitation control (5.5) based on the LS parameter estimate and (5.6) satisfies*

$$\limsup_{k \rightarrow \infty} \frac{1}{k} \sum_{i=1}^k \left(\|u_i\|^2 + \|y_i\|^2 \right) < \infty \quad \text{a.s. and} \quad \lim_{k \rightarrow \infty} \frac{1}{k} \sum_{i=1}^k (y_i - y_i^*)^2 = R \quad \text{a.s.,}$$

586 and the regressor φ_k satisfies the following excitation:

$$587 \quad (5.7) \quad \lambda_{\max}(n) \triangleq \lambda_{\max} \left\{ \sum_{k=1}^n \varphi_k \varphi_k^T \right\} = O(n), \quad \lambda_{\min}(n) \triangleq \lambda_{\min} \left\{ \sum_{k=1}^n \varphi_k \varphi_k^T \right\} \geq cn^{1-\bar{\varepsilon}(t+1)} \quad \text{a.s.}$$

588 for some $c > 0$, which may depend on sample paths and the $\bar{\varepsilon}$ defined below (5.6).

589 For the input and output signals generated by the system (5.4), by minimizing
 590 (3.2) in Algorithm 3.1, we can obtain the estimate of the sparse system parameters
 591 in (5.4). Denote the estimate as $\beta_{L,n} = [\beta_{L,n}(1), \dots, \beta_{L,n}(n_y + n_u)]^T$, and set

$$592 \quad H^* = \{i : a_i = 0 \text{ for } 1 \leq i \leq n_y; \text{ and } b_{i-n_y} = 0 \text{ for } n_y + 1 \leq i \leq n_y + n_u\},$$

$$593 \quad H_n^* = \{i : \beta_{L,n}(i) = 0 \text{ for } 1 \leq i \leq n_y + n_u\}.$$

594 Then, for the estimate $\beta_{L,n}$ obtained by Algorithm 3.1 with data $\{y_{k+1}, \varphi_k\}_{k=1}^n$ gen-
 595 erated from (5.4)-(5.6), the following theorem demonstrates the set convergence of
 596 the estimate in probability.

597 THEOREM 5.4. *If Assumptions (A1) and (D1)-(D3) hold, then*

$$598 \quad (5.8) \quad \lim_{n \rightarrow \infty} P(H_n^* = H^*) = 1,$$

599 where $\lambda_n = n^\tau$ in Algorithm 3.1 with $\tau \in \left(\frac{1}{2}\gamma + \frac{(1-\gamma)(2-\gamma)}{8-2\gamma}, \frac{1}{2}\right)$ and $\bar{\varepsilon} = \frac{1-\gamma}{8-2\gamma} \frac{1}{\bar{n}_{yu}+1}$ in
 600 the controller (5.6).

Proof. First, τ is well-defined, which can be verified by the following inequality:

$$\frac{1}{2} - \left(\frac{1}{2}\gamma + \frac{(1-\gamma)(2-\gamma)}{8-2\gamma} \right) = \frac{1}{2}(1-\gamma) - \frac{(1-\gamma)(2-\gamma)}{8-2\gamma} = (1-\gamma) \frac{1}{4-\gamma} > 0.$$

Denote $\lambda_{E,\max}(n) = \lambda_{\max} \left\{ E \sum_{k=1}^n \varphi_k \varphi_k^T \right\}$. From (5.7) in Proposition 5.3, we have
 $\lambda_{E,\max}(n) = O(n)$. Moreover, we can choose $d_n = c_1 n^{1-\bar{\varepsilon}(t+1)}$ with $c_1 \leq c$. By the

specification of $\bar{\varepsilon}$ below (5.6) and noting that $\tau < \frac{1}{2}$, we have $0 < \tau < \frac{1}{2} < 1 - \bar{\varepsilon}(t+1) < 1$. Then, it follows

$$\begin{aligned} \frac{\log \lambda_{\max}(n)}{\lambda_{\min}(n)} &= O\left(\frac{\log n}{n^{1-\bar{\varepsilon}(t+1)}}\right) \xrightarrow{n \rightarrow \infty} 0, \quad \frac{\lambda_n}{\lambda_{\min}(n)} = O\left(\frac{1}{n^{1-\bar{\varepsilon}(t+1)-\tau}}\right) \xrightarrow{n \rightarrow \infty} 0, \\ \frac{\lambda_n}{\lambda_{E,\max}(n)^{1/2}} &= O\left(\frac{1}{n^{\frac{1}{2}-\tau}}\right) \xrightarrow{n \rightarrow \infty} 0, \quad \frac{\sqrt{\lambda_{E,\max}(n)}}{d_n} = O\left(\frac{1}{n^{(1-\bar{\varepsilon}(t+1))-1/2}}\right) \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

Moreover, by noting $\bar{\varepsilon} = \frac{1-\gamma}{8-2\gamma} \frac{1}{t+1}$ and $\tau > \frac{1}{2}\gamma + \frac{(1-\gamma)(2-\gamma)}{8-2\gamma}$, we have $0 < \tau - \left(\frac{1}{2}\gamma + \frac{(1-\gamma)(2-\gamma)}{8-2\gamma}\right) = \tau - \left(\frac{1}{2}\gamma + \bar{\varepsilon}(t+1)(2-\gamma)\right)$, which implies

$$\begin{aligned} \frac{\lambda_n d_n^{2-\gamma}}{\lambda_{E,\max}(n)^{2-\frac{1}{2}\gamma}} &= O\left(n^{\tau+(1-\bar{\varepsilon}(t+1))(2-\gamma)-(2-\frac{1}{2}\gamma)}\right) \\ &= O\left(n^{\tau-(\frac{1}{2}\gamma+\bar{\varepsilon}(t+1)(2-\gamma))}\right) \xrightarrow{n \rightarrow \infty} \infty. \end{aligned}$$

By applying Theorem 3.10, the conclusion holds. \square

Remark 5.5. The weighted L_γ regularization Algorithm 4.1 can also be applied to these two typical problems, and the analyses are similar, and hence, omitted here.

6. Simulation study. This section sets up four simulations to validate the sparse identification performance of the proposed algorithms in this paper, including two finite impulse response (FIR) systems, a polynomial expansion NARX system and a linear feedback control system. In this paper, we use the particle swarm algorithm to solve (3.3) and (4.1).

Example 1. For the simulation of sparsity and estimation performance, consider the following FIR system: $y_{k+1} = \theta^T \varphi_k + w_{k+1}$, where $\theta = (1_{q \times 1}, 0_{(30-q) \times 1})^T$ with $q = 5, 10, 15, 20, 25$, φ_k are randomly generated in the interval $[-5, 5]$, and the noise sequence $\{w_k\}$ is i.i.d. with the Gaussian distribution $N(0, 0.1)$ and independent of $\{\varphi_k\}$. From Fig. 2, it can be seen that as the number of non-zero parameters q increases, the estimation error will be larger for the same number of samples, which also indicates that the smaller q is, the better the algorithm performs.

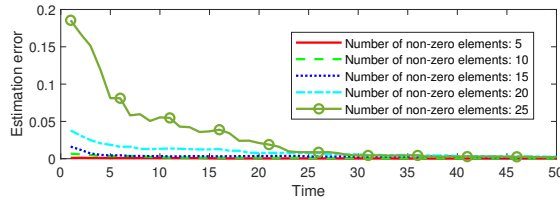


FIG. 2. Estimation error with different number of non-zero elements

Example 2. For the system (5.1), a common type of function expansion is the polynomial expansion [8], whose basis function is:

$$(6.1) \quad \varphi_j(x(k)) = y_{k-d_{j1}} \times \cdots \times y_{k-d_{ji}} \times u_{k-d_{j,i+1}} \times \cdots \times u_{k-d_{jl}},$$

where $d_{j1}, \dots, d_{jl} \in \mathbb{N}_+$, $l = 1, \dots, M$ with M being the maximum order of the polynomial expansion. Consider such a polynomial expansion NARX model, where $M = n_u = n_y = 2$. Then, the regressor $\varphi(k)$ contains $\frac{(M+n_y+n_u)!}{M!(n_y+n_u)!} = 15$ of possible basis functions. Let the real system be $y_{k+1} = \theta^T \varphi_k + w_{k+1}$, where $\theta = [\theta_1, \dots, \theta_{15}]$,

$$\begin{aligned} \varphi_k &= [u_k^2, u_k u_{k-1}, u_k y_k, u_k y_{k-1}, u_k, u_{k-1}^2, u_{k-1} y_k, \\ &\quad u_{k-1} y_{k-1}, u_{k-1}, y_k^2, y_k y_{k-1}, y_k, y_{k-1}^2, y_{k-1}, 1]^T \end{aligned}$$

TABLE 1

The objective function and parameter settings corresponding to the algorithms

Algorithm	Objective function	Algorithm parameters
LS	$\lambda_n = 0$	
LASSO	$\gamma = 1, \rho = 1$	$\lambda_n = n^{0.25}$
Ridge regression	$\rho = 0$	$\lambda_n = n^{0.05}$
Elastic net	$\gamma = 1$	$\rho = 0.5, \lambda_n = n^{0.25}$
Algorithm 3.1	$\rho = 1$	$\gamma = 0.4, \lambda_n = n^{0.25}$

619 The real parameters are set as $\theta = [0, -0.5, 0.7, 0, 0.45, 0, 0, -0.006, -0.5, 0, 0.008,$
620 $-0.2, 0, 1, 0]^T$. In this example, we use LS method ([6]), LASSO method ([33]), ridge
621 regression method ([16]), elastic net method ([44]), and Algorithm 3.1 to identify the
622 system parameters, respectively. A unified objective function of these methods takes
623 the following form

$$624 \quad (6.2) \quad J_{n+1}(\beta) = \sum_{k=1}^n (y_{k+1} - \beta^T \varphi_k)^2 + \lambda_n \rho \sum_{l=1}^q |\beta(l)|^\gamma + \lambda_n \frac{1-\rho}{2} \sum_{l=1}^q |\beta(l)|^2.$$

625 The form of the objective function corresponding to the algorithm and the param-
626 eter settings are given in Table 1. For this system, set the initial value to be i.i.d
627 with the input $\{u_k\}$, obeying the uniform distribution $U(-1, 1)$ and the noise $\{w_k\}$,
628 independent of $\{u_k\}$, obeying the normal distribution $N(0, 0.1)$.

629 Table 2 and Fig. 3 show the parameter estimation results of Algorithm 3.1, LS,
630 LASSO, ridge regression, and elastic net with 200 observations, respectively. From
631 Table 2 and Fig. 3, we can see that the Algorithm 3.1 has about the same accuracy
632 in estimating the non-zero parameters as the rest of the algorithms, but at the same
633 time, can significantly increase the accuracy of the selection of the zero parameters.
634 When $n = 200$, the approximation solution of the estimates of the zero parameters
635 are all less than 10^{-16} , indicating that Algorithm 3.1 performs better than the other
636 algorithms in identifying the zero parameters. Table 2 also shows the running time of
637 different methods. It is worth pointing out that the non-convex criterion adopted in
638 this paper greatly improves the identification accuracy although it inevitably increases
639 the computational complexity and the running time is relatively long.

TABLE 2

Comparison between Algorithm 3.1, LS, LASSO, Ridge regression and Elastic net under 200 observations.

Algorithms	$\theta_1 = 0$	$\theta_2 = -0.5$	$\theta_5 = 0.45$	$\theta_7 = 0$	Time
Algorithm 2.1	-2.2404×10^{-16}	-0.4950	0.4472	3.0363×10^{-17}	6.9296s
LS	-0.0011	-0.5010	0.4517	-0.0036	0.0228s
LASSO	-0.0015	-0.4958	0.4481	-7.6290×10^{-4}	0.5586s
Ridge regression	-6.9156×10^{-4}	-0.2805	0.3229	-0.0015	0.0348s
Elastic net	-0.0016	-0.4975	0.4493	-0.0022	0.6818s

640 **Example 3.** This example shows the application of the Algorithm 3.1 to the
641 identification of a linear feedback control system. Let the ARX system be

$$642 \quad (6.3) \quad y_{k+1} = \theta^T \varphi_k + w_{k+1} = \theta_1 y_k + \dots + \theta_5 y_{k+1-5} + \theta_6 u_k + \dots + \theta_{10} u_{k+1-5} + w_{k+1},$$

643 where the true sparse parameters are $\theta = [0.5, 3, 0, -1, 0.5, 0, 0, 0, 0, 0]^T$. The noise
644 $\{w_k\}$ is i.i.d, obeying the normal distribution $N(0, 0.025)$. The discrete reference

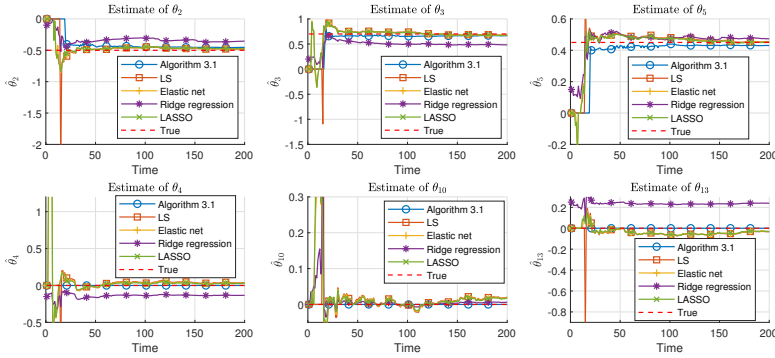


FIG. 3. Comparison between Algorithm 3.1, LS, LASSO, Ridge regression and Elastic net.

TABLE 3
Comparisons between Algorithm 3.1, LS, LASSO, Ridge regression and Elastic net at the 200th iteration.

	N=200		
Algorithm	$\theta_2 = 0$	$\theta_4 = 0$	$\theta_{10} = 0$
Algorithm 1	1.0433×10^{-17}	7.3991×10^{-17}	2.7050×10^{-17}
LS	0.0367	-0.0469	-0.1366
LASSO	0.0486	-0.0912	-0.1460
Ridge regression	0.0358	-0.0160	-0.0386
Elastic net	0.0481	-0.0210	-0.1476

645 signal is written as $y_{k+1}^* = \sin(\frac{1}{200}k)$, $k \geq 0$. Let the LS estimate be $\theta_k = [\theta_k(1), \dots,$
 646 $\theta_k(10)]^T$, then the self-tuning regulation control with diminishing excitation is

647 (6.4)
$$u_k = \frac{1}{\theta_k(6)} (y_{k+1}^* - (\theta_k(6)u_k - \theta_k^T \varphi_k)) + \frac{w'_k}{r_{k-1}^{\bar{\varepsilon}/2}},$$

648 where $r_{k-1} = 1 + \sum_{l=1}^{k-1} \|\varphi_l\|^2$, $\bar{\varepsilon} = \frac{1}{20}$ and $\{w'_k\}$ are i.i.d with the uniform distribution
 650 $U(-0.1, 0.1)$. Fig. 4 plots the outputs of the closed-loop control system (6.3)-(6.4)
 651 and the reference signals.

652 For the identification problem of the closed-loop control system, Table 3 and
 653 Fig. 5 show that, as long as excitation conditions are satisfied, Algorithm 3.1 can
 654 accurately distinguish between zero and non-zero parameters, and has more precise
 655 estimates of the zero parameters than other algorithms.

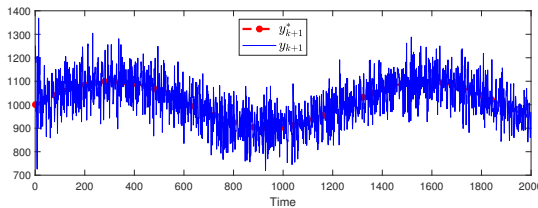


FIG. 4. Trajectories of y_{k+1} v.s. y_{k+1}^* for Example 2.

656 **Example 4.** This example aims to compare the performance of LS in [6], adaptive
 657 LASSO in [42] with Algorithm 4.1 in this paper. Consider the following FIR system:

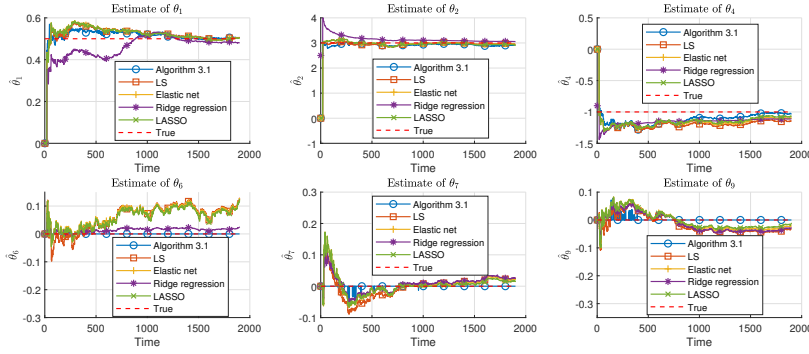


FIG. 5. Comparisons between Algorithm 3.1, LS, LASSO, Ridge regression and Elastic net.

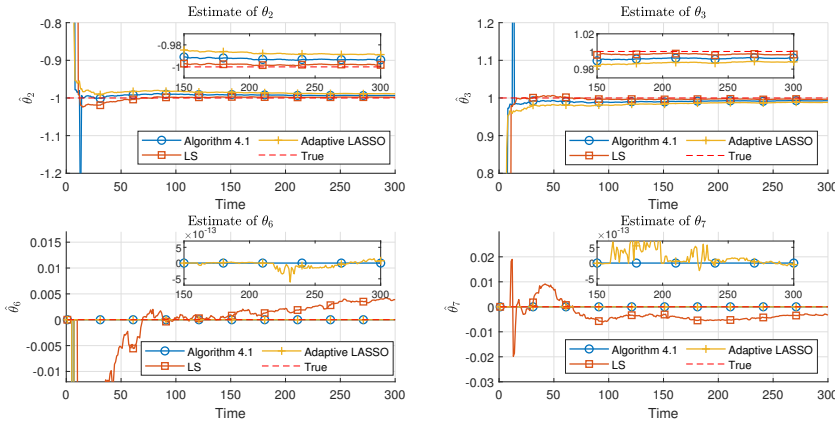


FIG. 6. Comparison between Algorithm 4.1, LS and adaptive LASSO.

658 $y_{k+1} = \theta^T \varphi_k + w_{k+1}$, where $\theta = [0, -1, 1, 2, 0.5, 0, 0, 0]^T$, φ_k are randomly generated in
 659 the interval $[-5, 5]$, and the noise sequence $\{w_k\}$ is i.i.d. with the Gaussian distribution
 660 $N(0, 0.1)$ and independent of $\{\varphi_k\}$. Set $\lambda_n = n^{0.65}$ for the adaptive LASSO in [42]
 661 and Algorithm 4.1. It can be seen from Fig. 6, Algorithm 4.1 provides a more sparse
 662 estimate of the system parameters than LS and the algorithm in [42], and a more
 663 accurate estimate than the adaptive LASSO in [42].

664 **7. Conclusion.** This paper investigates two kinds of sparse identification algo-
 665 rithms based on the non-convex L_γ penalty for the stochastic systems with non-i.i.d
 666 and non-stationary observation sequences and non-i.i.d noise. First, a one-step sparse
 667 parameter identification algorithm is proposed based on the L_γ ($0 < \gamma < 1$) penalty
 668 and the residual sum of squares. The almost sure convergence, the set convergence
 669 in probability, and the asymptotic normality property of the estimates generated by
 670 the proposed algorithm are established. Moreover, to improve the performance of
 671 the L_γ regularization method, a two-step algorithm based on the adaptively weighted
 672 L_γ ($0 < \gamma \leq 1$) penalty is provided. Not only is the almost sure parameter conver-
 673 gence of the estimates established, but also the almost sure set convergence is achieved.

674 Compared with existing literature, the theoretical results of the algorithms in this pa-
 675 per are applicable to the stochastic sparse system with non-i.i.d and non-stationary
 676 observation sequences and non-i.i.d noise and the algorithms are more efficient in
 677 sparsity induction. Furthermore, these algorithms are successfully implemented in
 678 the structure selection of the NARX models and the sparse parameter identification
 679 of the linear feedback control systems.

680 In the future, since sparsity is often accompanied by high dimensionality, it is
 681 interesting to consider the identification of stochastic sparse systems in high dimen-
 682 sional settings, i.e., $p = p(n)$. Moreover, it is essential to propose a recursive algorithm
 683 for the sparse system identification, and consequently, to design controls.

684

REFERENCES

- 685 [1] H. AKAIKE, *Information theory and an extension of the maximum likelihood principle*, Selected
 686 papers of hirotugu akaike, (1998), pp. 199–213.
- 687 [2] K. J. ÅSTRÖM AND B. WITTENMARK, *On self tuning regulators*, Automatica, 9 (1973), pp. 185–
 688 199.
- 689 [3] E. J. CANDÈS AND M. B. WAKIN, *An introduction to compressive sampling*, IEEE Signal Pro-
 690 cessing Magazine, 25 (2008), pp. 21–30.
- 691 [4] R. CHARTRAND AND V. STANEVA, *Restricted isometry properties and nonconvex compressive*
 692 *sensing*, Inverse Problems, 24 (2008), p. 035020.
- 693 [5] H. F. CHEN AND L. GUO, *Consistent estimation of the order of stochastic control systems*,
 694 IEEE Transactions on Automatic Control, 32 (1987), pp. 531–535.
- 695 [6] H. F. CHEN AND L. GUO, *Identification and stochastic adaptive control*, Springer Science &
 696 Business Media, 2012.
- 697 [7] A. DVORETZKY, *Asymptotic normality for sums of dependent random variables*, in Proceed-
 698 ings of the sixth Berkeley symposium on mathematical statistics and probability, vol. 2,
 699 University of California Press Berkeley, 1972, pp. 513–535.
- 700 [8] A. FALSONE, L. PIRODDI, AND M. PRANDINI, *A randomized algorithm for nonlinear model*
 701 *structure selection*, Automatica, 60 (2015), pp. 227–238.
- 702 [9] J. Q. FAN AND R. Z. LI, *Variable selection via nonconcave penalized likelihood and its oracle*
 703 *properties*, Journal of the American Statistical Association, 96 (2001), pp. 1348–1360.
- 704 [10] J. Q. FAN AND J. C. LV, *Nonconcave penalized likelihood with NP-dimensionality*, IEEE Trans-
 705 actions on Information Theory, 57 (2011), pp. 5467–5484.
- 706 [11] S. FOUCART AND M. J. LAI, *Sparsest solutions of underdetermined linear systems via l_q -*
 707 *minimization for $0 < q \leq 1$* , Applied and Computational Harmonic Analysis, 26 (2009),
 708 pp. 395–407.
- 709 [12] Y. X. FU AND W. X. ZHAO, *Support recovery and parameter identification of multivariate*
 710 *ARMA systems with Exogenous inputs*, SIAM Journal on Control and Optimization, 61
 711 (2023), pp. 1835–1860.
- 712 [13] A. GOLDSMITH, *Wireless communications*, Cambridge university press, 2005.
- 713 [14] L. GUO AND H. F. CHEN, *The Astrom-Wittenmark self-tuning regulator revisited and ELS-*
 714 *based adaptive trackers*, IEEE Transactions on Automatic Control, 36 (1991), pp. 802–812.
- 715 [15] L. GUO, H. F. CHEN, AND J. F. ZHANG, *Consistent order estimation for linear stochastic*
 716 *feedback control systems (CARMA model)*, Automatica, 25 (1989), pp. 147–151.
- 717 [16] A. E. HOERL AND R. W. KENNARD, *Ridge regression: Biased estimation for nonorthogonal*
 718 *problems*, Technometrics, 12 (1970), pp. 55–67.
- 719 [17] D. W. HUANG AND L. GUO, *Estimation of nonstationary ARMAX models based on the hannan-*
 720 *rissanen method*, The Annals of Statistics, 18 (1990), pp. 1729–1756.
- 721 [18] K. KNIGHT AND W. J. FU, *Asymptotics for LASSO-type estimators*, Annals of statistics, (2000),
 722 pp. 1356–1378.
- 723 [19] D. KRISHNAN AND R. FERGUS, *Fast image deconvolution using hyper-laplacian priors*, Advances
 724 in Neural Information Processing Systems, 22 (2009).
- 725 [20] M. J. LAI, Y. Y. XU, AND W. T. YIN, *Improved iteratively reweighted least squares for un-*
 726 *constrained smoothed l_q minimization*, SIAM Journal on Numerical Analysis, 51 (2013),
 727 pp. 927–957.
- 728 [21] T. L. LAI AND H. ROBBINS, *Consistency and asymptotic efficiency of slope estimates in sto-*
 729 *chastic approximation schemes*, Z. Wahrsch. verw. Gebiete, 56 (1981), pp. 329–360.
- 730 [22] T. L. LAI AND C. Z. WEI, *Least squares estimates in stochastic regression models with applica-*

- 731 *tions to identification and control of dynamic systems*, The Annals of Statistics, 10 (1982),
 732 pp. 154–166.
- 733 [23] T. L. LAI AND C. Z. WEI, *On the concept of excitation in least squares identification and adap-*
 734 *tive control*, Stochastics: An International Journal of Probability and Stochastic Processes,
 735 16 (1986), pp. 227–254.
- 736 [24] J. H. LIN AND G. MICHAILEDIS, *System identification of high-dimensional linear dynamical*
 737 *systems with serially correlated output noise components*, IEEE Transactions on Signal
 738 Processing, 68 (2020), pp. 5573–5587.
- 739 [25] N. LIU, W. LI, Y. J. WANG, R. TAO, Q. DU, AND J. CHANUSSOT, *A survey on hyperspec-*
 740 *tral image restoration: From the view of low-rank tensor approximation*, Science China
 741 Information Sciences, 66 (2023), pp. 1–31.
- 742 [26] L. LJUNG, *System identification*, Wiley encyclopedia of electrical and electronics engineering,
 743 (1999), pp. 1–19.
- 744 [27] K. LU, H. LIU, L. ZENG, J. Y. WANG, Z. S. ZHANG, AND J. P. AN, *Applications and prospects*
 745 *of artificial intelligence in covert satellite communication: a review*, Science China Infor-
 746 mation Sciences, 66 (2023), pp. 1–31.
- 747 [28] N. MEINSHAUSEN AND P. BÜHLMANN, *Variable selection and high-dimensional graphs with the*
 748 *lasso*, The Annals of Statistics, 34 (2006), pp. 1436–1462.
- 749 [29] J. K. PANT, W. S. LU, AND A. ANTONIOU, *New improved algorithms for compressive sensing*
 750 *based on l_p norm*, IEEE Transactions on Circuits and Systems II: Express Briefs, 61 (2014),
 751 pp. 198–202.
- 752 [30] M. M. PETROU AND C. PETROU, *Image processing: the fundamentals*, John Wiley & Sons,
 753 2010.
- 754 [31] A. ROSS AND A. JAIN, *Information fusion in biometrics*, Pattern Recognition Letters, 24 (2003),
 755 pp. 2115–2125.
- 756 [32] G. SCHWARZ, *Estimating the dimension of a model*, The Annals of Statistics, 6 (1978), pp. 461–
 757 464.
- 758 [33] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical
 759 Society: Series B (Methodological), 58 (1996), pp. 267–288.
- 760 [34] R. TÓTH, B. M. SANANDAJI, K. POOLLA, AND T. L. VINCENT, *Compressive system identifi-*
 761 *cation in the linear time-invariant framework*, in 2011 50th IEEE Conference on Decision
 762 and Control and European Control Conference, IEEE, 2011, pp. 783–790.
- 763 [35] A. WÄCHTER AND L. T. BIEGLER, *On the implementation of an interior-point filter line-search*
 764 *algorithm for large-scale nonlinear programming*, Mathematical Programming, 106 (2006),
 765 pp. 25–57.
- 766 [36] J. WOODWORTH AND R. CHARTRAND, *Compressed sensing recovery via nonconvex shrinkage*
 767 *penalties*, Inverse Problems, 32 (2016), p. 075004.
- 768 [37] Z. B. XU, X. Y. CHANG, F. M. XU, AND H. ZHANG, *$L_{1/2}$ regularization: A thresholding rep-*
 769 *resentation theory and a fast solver*, IEEE Transactions on Neural Networks and Learning
 770 Systems, 23 (2012), pp. 1013–1027.
- 771 [38] Z. B. XU, H. ZHANG, Y. WANG, X. Y. CHANG, AND Y. LIANG, *$L_{1/2}$ regularization*, Science
 772 China Information Sciences, 53 (2010), pp. 1159–1169.
- 773 [39] L. T. ZHANG AND L. GUO, *Adaptive identification with guaranteed performance under saturated*
 774 *observation and nonpersistent excitation*, IEEE Transactions on Automatic Control, 69
 775 (2024), pp. 1584–1599.
- 776 [40] P. ZHAO AND B. YU, *On model selection consistency of Lasso*, The Journal of Machine Learning
 777 Research, 7 (2006), pp. 2541–2563.
- 778 [41] W. X. ZHAO, *Parametric identification of Hammerstein systems with consistency results using*
 779 *stochastic inputs*, IEEE Transactions on Automatic Control, 55 (2010), pp. 474–480.
- 780 [42] W. X. ZHAO, G. YIN, AND E.-W. BAI, *Sparse system identification for stochastic systems with*
 781 *general observation sequences*, Automatica, 121 (2020), p. 109162.
- 782 [43] H. ZOU, *The adaptive LASSO and its oracle properties*, Journal of the American Statistical
 783 Association, 101 (2006), pp. 1418–1429.
- 784 [44] H. ZOU AND T. HASTIE, *Regularization and variable selection via the elastic net*, Journal of
 785 the Royal Statistical Society: series B (statistical methodology), 67 (2005), pp. 301–320.