

Regularized Nonlinear Regression for Simultaneously Selecting and Estimating Key Model Parameters

Kyubaek Yoon · Hojun You · Wei-Ying Wu · Chae Young Lim · Jongeun Choi ·
Connor Boss · Ahmed Ramadan · John M. Popovich Jr. · Jacek Cholewicki ·
N. Peter Reeves · Clark J. Radcliffe

Received: date / Accepted: date

Abstract In system identification, estimating parameters of a model using limited observations results in poor identifiability. To cope with this issue, we propose a new method to simultaneously select and estimate sensitive parameters as key model parameters and fix the remaining parameters to a set of typical values. Our method is formulated as a nonlinear least squares estimator with L_1 -regularization on the deviation of parameters from a set of typical values. First, we provide consistency and oracle properties of the proposed estimator as a theoretical foundation. Second, we provide a novel approach based on Levenberg-Marquardt opti-

mization to numerically find the solution to the formulated problem. Third, to show the effectiveness, we present an application identifying a biomechanical parametric model of a head position tracking task for 10 human subjects from limited data. In a simulation study, we analyze the bias and variance of estimated parameters. In an experimental study, our method improves the model interpretation by reducing the number of parameters to be estimated while maintaining variance accounted for (VAF) at above 82.5%. Moreover, the variance of estimated parameters is reduced by 71.1% as compared to that of the estimated parameters without L_1 -regularization. Our method is 54 times faster than the standard simplex-based optimization to solve the regularized nonlinear regression.

Keywords System identification · Nonlinear regression · L_1 -regularization · Lasso · Levenberg-Marquardt Optimization

1 Introduction

In parameter estimation, a model is considered to be identifiable when a unique set of parameters is specified for given measurement data. However, when the data is limited, estimating unknown parameters of a model results in poor identifiability [7]. In such a case, small changes in the data could result in very different estimated parameters, for example, a rather randomly chosen local minimum out of multiple local minima [14, 20]. The resulted overfitting impairs the model parsimony and generalizability [17]. The overfitted model may yield good results with a training data set used to estimate parameters, but it may yield poor estimates with a new test data set. Moreover, this issue becomes worse when parameters are estimated from the data corrupted by random noise [?].

Kyubaek Yoon · Jongeun Choi (Corresponding author)
The School of Mechanical Engineering, Yonsei University, 50 Yonsei
Ro, Seodaemun Gu, Seoul 03722, Republic of Korea
E-mail: jongeunchoi@yonsei.ac.kr

Hojun You · Chae Young Lim
Department of Statistics, Seoul National University, Seoul 08826, Republic of Korea

Wei-Ying Wu
Department of Applied Mathematics, National Dong Hwa University,
Hualien 97401, Taiwan

Connor Boss
Department of Electrical Engineering, Michigan State University, East
Lansing, MI 48824, USA

Ahmed Ramadan
Department of Physical Therapy and Rehabilitation Science, University
of Maryland, Baltimore, MD 21201, USA

John M. Popovich Jr. · Jacek Cholewicki
MSU Center for Orthopedic Research, Department of Osteopathic Surgical
Specialties, Michigan State University, East Lansing, MI 48824
USA

N. Peter Reeves
Sumaq Life LLC, East Lansing, MI 48823 USA

Clark J. Radcliffe
Department of Mechanical Engineering and MSU Center for Orthopedic
Research, Michigan State University, East Lansing, MI 48824
USA

A biomechanical model often has unknown parameters to be estimated with limited data due to unavailable internal states and the non-invasive nature of human data collection [13]. For a limited observation data set, one way for improving identifiability is to build a parsimonious model by lumping a large number of parameters into a small number of lumped parameters [4]. Such a parsimonious model has better interpretability and provides higher estimation accuracy for arbitrary data [1].

As a way to build a parsimonious model, the Least absolute shrinkage and selection operator (Lasso) was first introduced in [23] and then further developed in [24, 27, 29, 30]. The Lasso is typically used to select sensitive parameters among parameters of a linear model. A sensitive parameter subset is considered to have relatively large impact on the output of the model. That is, small changes in the sensitive parameters result in large changes in the model response [25].

To formally introduce the Lasso, we suppose that we have $(\mathbf{x}_i, y_i), i \in \{1, \dots, n\}$ where $\mathbf{x}_i = [x_i^{[1]}, \dots, x_i^{[p]}]$ and $y_i = f(\mathbf{x}_i; \boldsymbol{\theta}_0) + \epsilon_i$. $f(\mathbf{x}_i; \boldsymbol{\theta}_0)$ is a function, which depends on the true parameter vector $\boldsymbol{\theta}_0$. ϵ_i is independent and identically distributed with $\mathbb{E}(\epsilon) = 0$ and $Var(\epsilon_i) = \sigma^2$. Without loss of generality, we assume that the true parameter vector $\boldsymbol{\theta}_0 = [\theta_{01}, \theta_{02}, \dots, \theta_{0s}, \theta_{0s+1}, \dots, \theta_{0p}]^T$ has the first s entries non-zero. That is, $\theta_{0k} \neq 0$, for $1 \leq k \leq s$ and $\theta_{0k} = 0$, for $s+1 \leq k \leq p$. Finally, when $f(\mathbf{x}_i; \boldsymbol{\theta}_0) = \mathbf{x}_i \boldsymbol{\theta}_0$, we consider the following linear least squares problem with L_1 -regularization.

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \left[\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 + n\lambda \sum_{k=1}^p |\theta_k| \right], \quad (1)$$

where $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{n \times p}$ is the input. $\mathbf{y} = [y_1, y_2, \dots, y_n]^T \in \mathbb{R}^{n \times 1}$ is the observation vector. In (1), the hyperparameter $\lambda > 0 \in \mathbb{R}$ determines the amount of regularization. The Lasso shrinks more number of parameters toward 0 as λ increases in general. Moreover, insensitive parameters are shrunk to 0 if λ is sufficiently large [29]. The remaining parameters, which are not shrunk to 0, are considered sensitive parameters [20]. In this paper, we consider the sensitive parameters as key model parameters. The L_1 -regularization methods and weighted least squares method were used to select nonlinear autoregressive with exogenous variables (NARX) models [19]. The Lasso was used to remove insensitive parameters when all unknown parameters are to be non-negative [11]. The modified Lasso was proposed for the nonlinear induction motor identification problem, which deals with a similar problem to our study [22]. However, they do not provide consistency and the asymptotic normality results. The modified principal component analysis (PCA) based on the Lasso was proposed for dimension reduction [10]. To select and estimate sensitive parameters of a model, sensitive parameters were selected based

on parameter estimation variances predicted by the Fisher information matrix [20, 21].

In this paper, our objective is to develop a regularized nonlinear parameter estimation method for a model with unknown parameters using a limited data set. Thus, we formulate the model parameter estimation problem as a nonlinear least squares problem with L_1 -regularization as follows.

$$\hat{\tilde{\boldsymbol{\theta}}} = \arg \min_{\tilde{\boldsymbol{\theta}}} \left[\|\mathbf{y} - f(\mathbf{X}; \tilde{\boldsymbol{\theta}})\|_2^2 + n\lambda \sum_{k=1}^p |\tilde{\theta}_k| \right], \quad (2)$$

where $\tilde{\boldsymbol{\theta}} := \boldsymbol{\theta} - \bar{\boldsymbol{\theta}}$ and $\bar{\boldsymbol{\theta}}$ is the set of typical values. Note that these values may be obtained as the mean values of $\boldsymbol{\theta}$ based on preliminary information or estimation. The details of (2) are introduced in Section 3. In our preliminary work, we developed a parameter selection method for system identification with application to head-neck position tracking and reported the model parameter estimates [12].

The contributions of the paper are as follows. First, we consider nonlinear regression with a generalized penalty function that includes an L_1 -penalty function and provide its consistency and oracle properties (i.e., convergence to the correct sparsity and asymptotic normality) in Section 2. To the best of the authors' knowledge, our work is the first to provide such analysis for nonlinear regression with a generalized penalty function. For example, convergence properties for various penalized linear regression have been discussed [6, 9]. Note that we do not assume the distribution of errors, which is different from the assumption of [6]. In Section 3, we then reformulate the regularized nonlinear regression for simultaneously selecting and estimating key model parameters by defining $\tilde{\theta}_k$ in (2) as the deviation of the k -th parameter from its nominal value. Next, we improve the optimization algorithm of [22] to numerically solve the nonlinear least squares problem. Finally, to show the effectiveness, we present an application identifying a biomechanical parametric model of a head-neck position tracking task from limited data in simulation and experimental studies. In a simulation study, our algorithm reduces the variance of most parameter estimates as well as the bias. In an experimental study, the variance of selected sensitive parameters is reduced by 71.1% on average while maintaining the goodness of fit at above 82.5%. In addition, our method is 54 times faster as compared to the Lasso with the brute force optimization, e.g., the standard simplex-based optimization we presented recently [20].

2 Consistency and Oracle Properties

We propose a penalized nonlinear regression approach for parameter selection and estimation. First of all, we adopt the following equivalent nonlinear least squares estimator with

a generalized penalty function from (2):

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in D} \left[\mathbf{Q}_n(\boldsymbol{\theta}) := \mathbf{S}_n(\boldsymbol{\theta}) + n \sum_{k=1}^p p_{\lambda_n}(|\theta_k|) \right], \quad (3)$$

where $\mathbf{S}_n(\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2$. The first term in $\mathbf{Q}_n(\cdot)$ corresponds to nonlinear least squares estimation, and the second term, $p_{\lambda_n}(\cdot)$ is the penalty function used for parameter selection. λ_n in the penalty function is a nonnegative regularization parameter. Note that we consider a generalized penalty function in (3). The appropriate choice of the penalty function, including L_1 -regularization, is further investigated in the assumption 2.

Without loss of generality, we assume only a few parameters are non-zero, such that the true parameter vector $\boldsymbol{\theta}_0 = [\theta_{01}, \theta_{02}, \dots, \theta_{0s}, \theta_{0s+1}, \dots, \theta_{0p}]$ has the first s entries non-zero. That is, $\theta_{0k} \neq 0$, for $1 \leq k \leq s$ and $\theta_{0k} = 0$, for $s+1 \leq k \leq p$. For the nonlinear function $f(\cdot; \cdot)$ in (3), we further consider the following assumptions.

Assumption 1 (Nonlinear function)

1. The true parameter $\boldsymbol{\theta}_0$ is in the interior of the bounded parameter set Θ , and $f(\mathbf{x}_i; \boldsymbol{\theta})$ is twice differentiable with respect to $\boldsymbol{\theta}$ near $\boldsymbol{\theta}_0$ for all i .
2. Let $\mathbf{f}_k(\boldsymbol{\theta}) = \left(\frac{\partial f(\mathbf{x}_1, \boldsymbol{\theta})}{\partial \theta_k}, \dots, \frac{\partial f(\mathbf{x}_n, \boldsymbol{\theta})}{\partial \theta_k} \right)^T$ and $\dot{\mathbf{F}}(\boldsymbol{\theta}) = (\mathbf{f}_1, \dots, \mathbf{f}_p)$. Then, there exists a positive definite matrix Γ such that $\frac{1}{n} \dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \rightarrow \Gamma$ as $n \rightarrow \infty$.
3. As $n \rightarrow \infty$ and $\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_0\| \rightarrow 0$,

$$\dot{\mathbf{F}}(\boldsymbol{\theta}_1)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_1) \left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1} \rightarrow I$$

uniformly, where I is the identity matrix

4. There exists a $\delta > 0$ such that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \leq \delta} \left(\frac{\partial^2 f(\mathbf{x}_i; \boldsymbol{\theta})}{\partial \theta_k \partial \theta_s} \right)^2 < \infty$$

For the penalty function $p_{\lambda_n}(\cdot)$ in (3), we further consider the following assumptions.

Assumption 2 (Penalty function) The first and second derivative of the penalty function $p_{\lambda_n}(\cdot)$ denoted by $q_{\lambda_n}(\cdot)$ and $q'_{\lambda_n}(\cdot)$ have the following properties:

1. For a nonzero fixed θ ,

$$\lim_{n \rightarrow \infty} n^{1/2} q_{\lambda_n}(|\theta|) = 0, \quad \lim_{n \rightarrow \infty} q'_{\lambda_n}(|\theta|) = 0.$$

2. For any $M > 0$,

$$\lim_{n \rightarrow \infty} n^{1/2} \inf_{|\theta| \leq M n^{-1/2}} q_{\lambda_n}(|\theta|) \rightarrow \infty.$$

Remark 1 Assumption 2 is satisfied for several well known penalty functions, e.g., SCAD, Adaptive Lasso, and Hard penalty with proper choices of λ_n . Assumption 2-(1) is satisfied for L_1 -regularization with a proper choice of λ_n . The details are discussed in [9].

The following theorem shows the existence of a local minimizer of $\mathbf{Q}_n(\boldsymbol{\theta})$ with the order of $O_p(n^{-1/2})$.

Lemma 1 Under Assumptions 1 and 2-(1), for any $\eta > 0$, there exists a positive constant C that makes, for large enough n ,

$$P \left(\inf_{\|\mathbf{v}\|=C} \mathbf{Q}_n(\boldsymbol{\theta}_0 + n^{-1/2} \mathbf{v}) - \mathbf{Q}_n(\boldsymbol{\theta}_0) > 0 \right) > 1 - \eta,$$

where $\mathbf{Q}_n(\boldsymbol{\theta}) = \mathbf{S}_n(\boldsymbol{\theta}) + n \sum_{k=1}^p p_{\lambda_n}(|\theta_k|)$.

Theorem 1 Under the assumptions in Lemma 1, there exists, with probability tending to 1, a root- n -consistent local minimizer $\hat{\boldsymbol{\theta}}$ of $\mathbf{Q}_n(\boldsymbol{\theta})$, that is, $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| = O_p(n^{-1/2})$.

Next theorem shows oracle properties (i.e., convergence to the correct sparsity and asymptotic normality) of the estimator on the true set.

Theorem 2 Assume that $\hat{\boldsymbol{\theta}} = (\hat{\theta}_k)_{k=1}^p$ is the local minimizer of $\mathbf{Q}_n(\boldsymbol{\theta})$ with the root- n -consistency. If Assumptions 1 and 2 hold,

- (i) for the set $M_k = \{\omega : \hat{\theta}_k \neq 0\}$, $s+1 \leq k \leq p$,

$$P(M_k) \rightarrow 0$$

- (ii) For $\hat{\boldsymbol{\theta}}_{11} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_s)^T$, $\boldsymbol{\theta}_{01} = (\theta_{01}, \theta_{02}, \dots, \theta_{0s})^T$,

$$n^{1/2}(2\Gamma_{11})(\hat{\boldsymbol{\theta}}_{11} - \boldsymbol{\theta}_{01}) + (2\Gamma_{11})^{-1} b_n \xrightarrow{d} N(0, 2\Gamma_{11}\sigma^2),$$

where $b_n = (q_{\lambda_n}(|\theta_{01}|) \text{sgn}(\theta_{01}), q_{\lambda_n}(|\theta_{02}|) \text{sgn}(\theta_{02}), \dots, q_{\lambda_n}(|\theta_{0s}|) \text{sgn}(\theta_{0s}))^T$ and Γ_{11} is the first $s \times s$ submatrix of Γ .

The proofs of Lemma 1 and Theorem 1 are given in Appendix A and Appendix B, respectively and the proof of Theorem 2 is given in Appendix C.

3 Application to Head Position Tracking

In this section, we evaluate our approach to solve the nonlinear least squares problem with L_1 -regularization from simulation and experimental studies of a biomechanical parametric model of a head-neck position tracking task in [20]. The reliability of the head-neck position tracking task to quantify head-neck motor control is demonstrated in [18]. As compared to the Levenberg optimization algorithm in [22], our method implements the Levenberg-Marquardt optimization algorithm to numerically solve our nonlinear least squares problem with L_1 -regularization. The Levenberg-Marquardt optimization uses the diagonal elements of the hessian matrix approximation to overcome the slow convergence problem when the value of the damping factor is large [15]. The

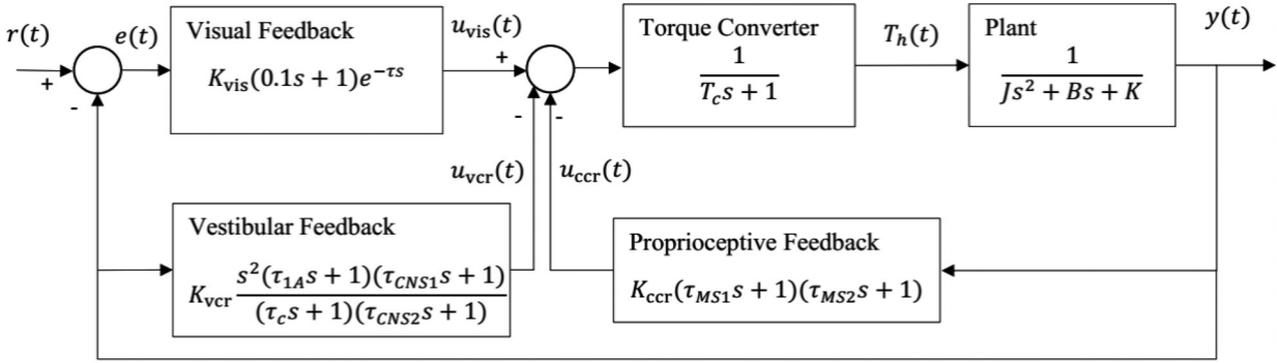


Fig. 1 Sensorimotor control block diagram for the head-neck system [20]

details of our algorithm are shown in Appendix D and Table 3.

Additionally, in order to simultaneously select and estimate key model parameters of the head-neck system, we reformulate the Lasso penalty function as L_1 -regularization on the deviation of parameters from a set of typical values. In this paper, we adopt the mean of the nominal parameter values obtained from preliminary estimation as the set of typical values. Therefore, our method simultaneously selects and estimates only sensitive parameters while fixing insensitive parameters onto the mean of the nominal parameter values. In order to compare with the method of [20], we set the number of sensitive parameters to 5. In this case, we increase the regularization hyperparameter value until we obtain 5 sensitive parameters. The goodness of fit is quantitatively evaluated by variance accounted for (VAF). VAF represents how much the experimental data can be explained by a fitted regression model. VAF is formally defined as follows.

$$\text{VAF}(\boldsymbol{\theta})(\%) = \left[1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i(\boldsymbol{\theta}))^2}{\sum_{i=1}^n y_i^2} \right] \times 100$$

3.1 Subjects

10 healthy subjects participated in the experimental study. They did not have any history of neck pain lasting more than three days or any neurological motor control impairments. The Michigan State University's Biomedical and Health Institutional Review Board approved the test protocol. The subjects signed an informed consent before participating in the experiment [20, 21].

3.2 Parametric model

Fig. 1 shows the block diagram of the head-neck system for position tracking. This is a representative physiological

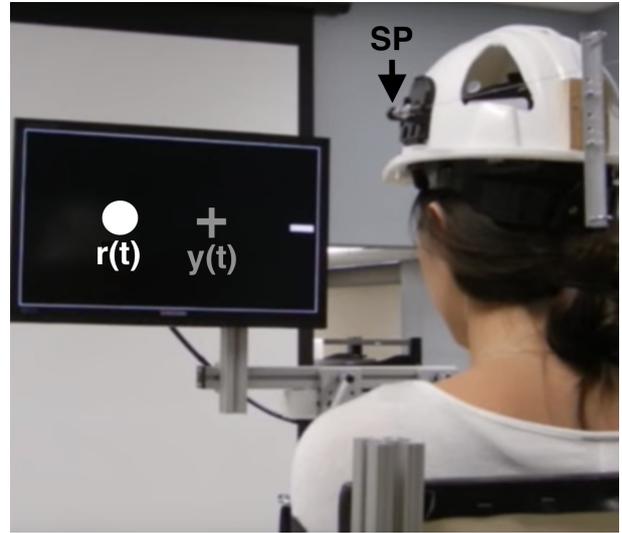


Fig. 2 The experimental setup for the head-neck position tracking task. $r(t)$ is the reference command signal, and $y(t)$ is the measured head rotation angle, and SP indicates one of string potentiometers on both sides of the helmet.

feedback control model [2, 16]. The model consists of 14 parameters. As shown in Table 1, 2 out of 14 parameters are set to fixed values from [16]. The remaining 12 parameters to be estimated are:

$$\boldsymbol{\theta} = [K_{vis} \ K_{vcr} \ K_{ccr} \ \tau \ \tau_{1A} \ \tau_{CNS1} \ \tau_C \ \tau_{CNS2} \ \tau_{MS1} \ \tau_{MS2} \ B \ K]$$

The remaining 12 parameters have the lower and upper bounds from [20] and are normalized using min-max normalization in order to ignore the scale differences between parameters.

3.3 The Experiment

As shown in Fig. 2, each subject wears a helmet attached with two string potentiometers measuring the axial rotation

Table 1 The neurophysiological parameters of the head position model. All the information is obtained from [20]

Parameters	Max	Min	Description
$K_{vis}[\frac{Nm}{rad}]$	10^3	50	Visual feedback gain
$K_{vcr}[\frac{Nm s^2}{rad}]$	10^4	500	Vestibular feedback gain
$K_{ccr}[\frac{Nm}{rad}]$	300	1	Proprioceptive feedback gain
$\tau[s]$	0.4	0.1	Visual feedback delay
$\tau_{1A}[s]$	0.2	0.01	Lead time constant of the irregular vestibular afferent neurons
θ $\tau_{CNS1}[s]$	1	0.05	Lead time constant of the central nervous system
$\tau_C[s]$	5	0.1	Lag time constant of the irregular vestibular afferent neurons
$\tau_{CNS2}[s]$	60	5	Lag time constant of the central nervous system
$\tau_{MS1}[s]$	1	0.01	First lead time constant of the neck muscle spindle
$\tau_{MS2}[s]$	1	0.01	Second lead time constant of the neck muscle spindle
$B[\frac{Nm}{rad}]$	5	0.1	Intrinsic damping
$K[\frac{Nm}{rad}]$	5	0.1	Intrinsic stiffness
θ_{fixed} $J[kgm^2]$	0.0148	0.0148	Head inertia
$T_c[s]$	0.1	0.1	Torque converter time constant

Subject	Improvement (%)
1	69.25
2	65.79
3	74.83
4	89.36
5	87.44
6	65.02
7	53.62
8	77.70
9	63.44
10	67.62
Avg.	71.41

Table 2 Improvement on the variance of estimated parameters in an experiment study. The given values are $(1 - \frac{\sigma^{2,Lasso}}{\sigma^{2,All}}) \times 100$. $\sigma^{2,Lasso}$ is the variance of estimated parameters obtained from our method, and $\sigma^{2,All}$ is the variance of estimated parameters without regularization.

of the head. Subjects rotated their heads about the vertical axis to track the command signals on the display. The command signal $r(t)$ is a pseudorandom sequence of steps with random step durations and amplitudes. The angle of the signal is bounded between $\pm 4^\circ$. The output signal $y(t)$ is the head rotation angle. Each subject performed three 30-second trials, and the sampling rate was 60 Hz [20].

3.4 Simulation study

In this section, we analyze the bias and variance of estimated parameters from a simulation study with the known true parameters for comparison. First, we generate the simulated data (\mathbf{x}, \mathbf{y}) where $\mathbf{x} \in \mathbb{R}^{1800 \times 1}$ and $\mathbf{y} \in \mathbb{R}^{1800 \times 1}$, which are the input and observation vectors, respectively. Additionally, we obtain 20 sets of nominal parameter values from preliminary estimation performed 20 times over all three trials for each subject. Finally, we evaluate our method by setting the mean of 20 sets of nominal parameter values as a set of typical values. In Fig. 3, as compared to the nonlinear least squares estimator without L_1 -regularization, except for τ_c ,

the biases of the other parameters were decreased by 28.0% on average. In addition, the variances of all estimated parameters were decreased by 96.1% on average.

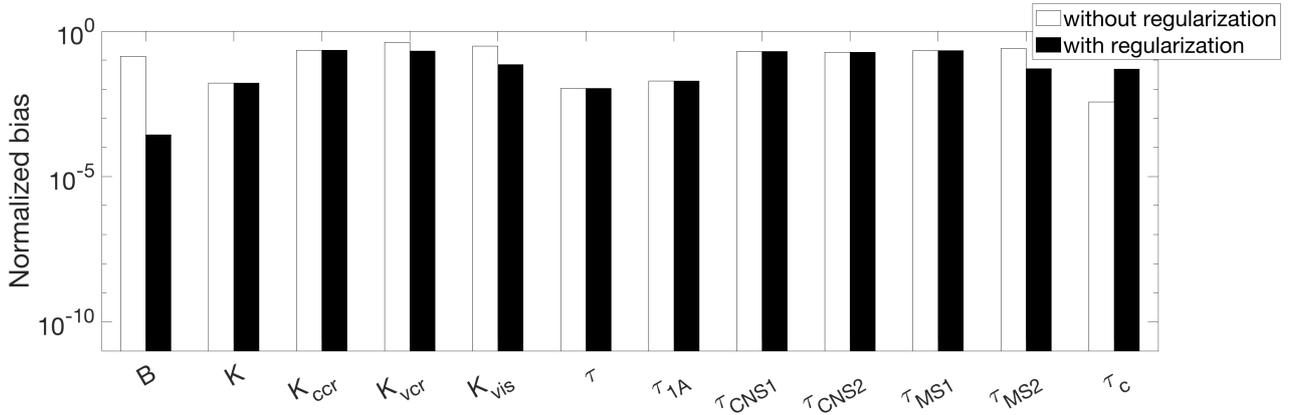
3.5 In vivo experimental study

In this section, the variances of estimated parameters from our method are compared with the result of the standard simplex-based optimization in [20]. All parameters are pushed further toward the mean of nominal parameter values obtained from preliminary estimation as the regularization hyperparameter increases. The regularization hyperparameter increases until 5 sensitive parameters are selected. After selecting a sensitive parameter subset for each subject, we select the most frequent subset among all subjects for the fair comparison with [20]. The Lasso in [20] selected 5 parameters of $[K_{vcr} K_{ccr} \tau_{1A} \tau_c \tau_{CNS2}]$ as the most frequent subset of sensitive parameters, and the subset of $[K_{ccr} \tau \tau_{1A} \tau_c \tau_{CNS2}]$ was selected by our method. As a result, 4 out of 5 sensitive parameters $[K_{ccr} \tau_{1A} \tau_c \tau_{CNS2}]$ are selected by both methods.

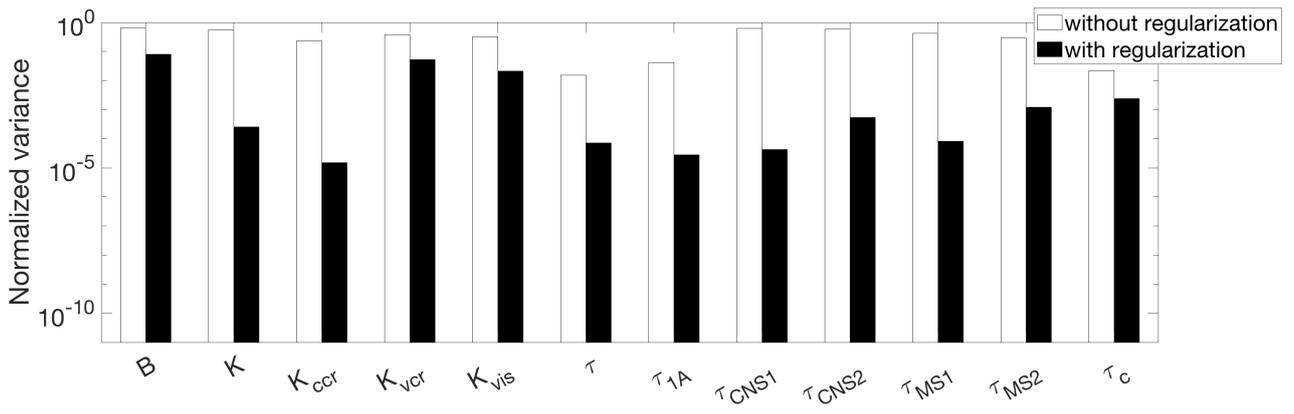
Second, we evaluate our method based on the goodness of fit measured by variance accounted for (VAF). All values are given as mean \pm standard deviation across subjects. The goodness of fit of our method (VAF = $82.5\% \pm 8.3\%$) over 10 subjects is almost equal to that of the Lasso (VAF = $83.3\% \pm 7.3\%$) in [20]. Without L_1 -regularization, VAF = $84.9\% \pm 0.1\%$ over all subjects. In Fig. 4, for all subjects, the estimated responses are almost same as the measured responses, and the estimated responses are smoother than the measured responses.

Third, as shown in Table 2, the variance of estimated parameters is reduced by 71.4% on average across parameters and subjects as compared to those of estimated parameters without regularization.

Finally, we compute the average computation time using the “timeit” function from *MATLAB* (The MathWorks



(a)



(b)

Fig. 3 The bias (a) and variance (b) of estimated parameters in the simulation studies. Comparing our method with a nonlinear least squares problem without L_1 -regularization (white bar), we set the mean of nominal parameters (black bar) obtained from preliminary estimation as a set of typical values. The y-axis is a logarithmic scale.

Inc., Natick, MA, U.S.A.). The average computation time of our method for a subject is 54 times faster than that of the Lasso with the standard simplex-based optimization in [20]. In particular, the average computation time of our method across subjects is 5.6 seconds per trial, and that of the Lasso in [20] is 302.0 seconds per trial.

4 Discussion

We provided consistency and oracle properties (i.e., convergence to the correct sparsity and asymptotic normality) for a nonlinear regression approach with a generalized penalty

function. As a result, we proved the existence of a local minimizer and the convergence to the sparse unknown parameters for the penalized nonlinear least squares estimator. It is important to note that for the first time, we have proved convergence properties of the penalized nonlinear least squares estimator, as compared to previous studies [6, 9, 26].

In the simulation study, we showed that the bias and variance of estimated parameters of our method were decreased as compared to those of estimated parameters without L_1 -regularization. When we set the mean of nominal parameter values as the typical values for non-selected estimates, the variance significantly was decreased. In addition, although

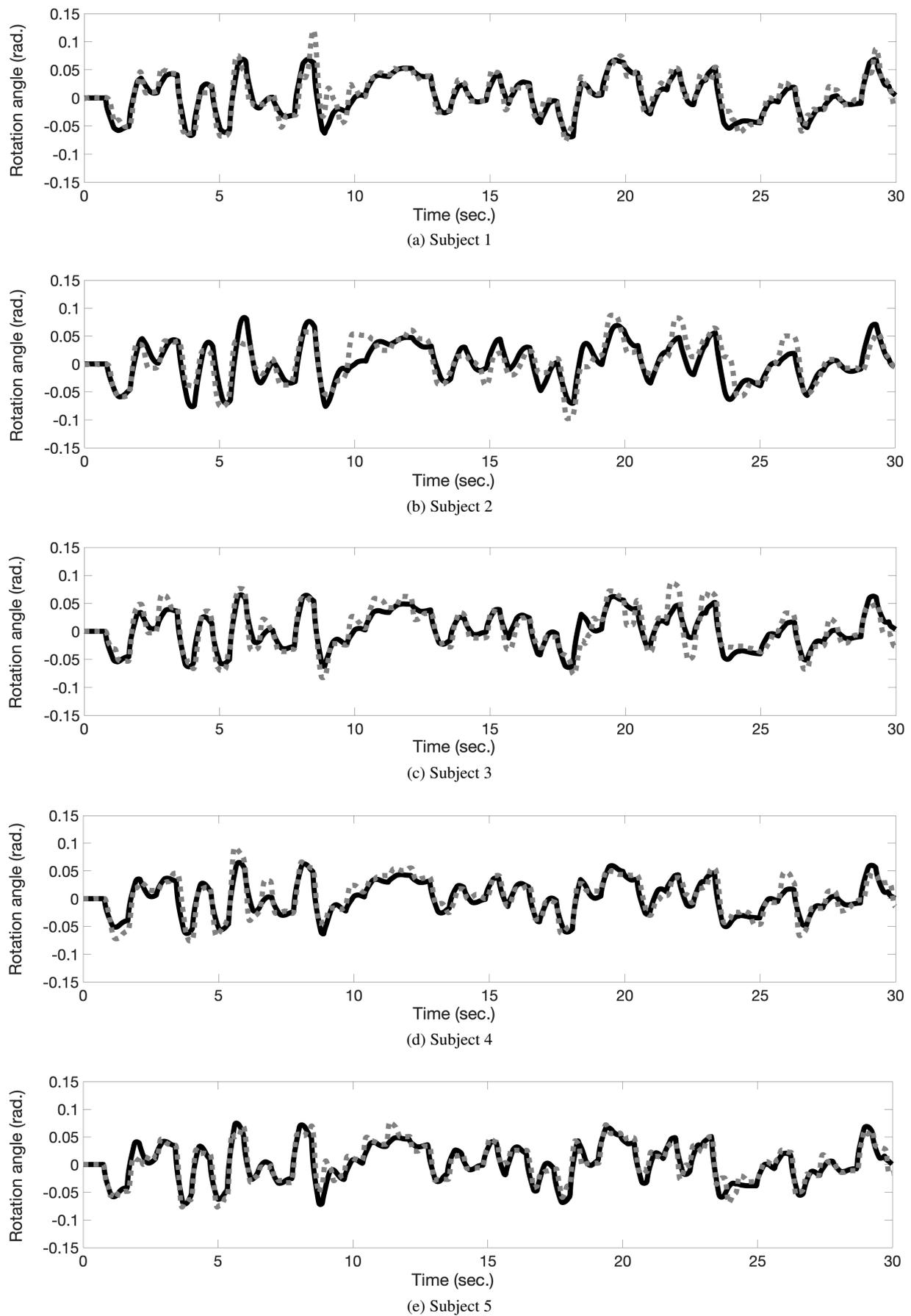


Fig. 4 The curve fitting with 10 experimental cases. Solid lines represent estimated responses from the fitted models and dotted lines represent measured responses.

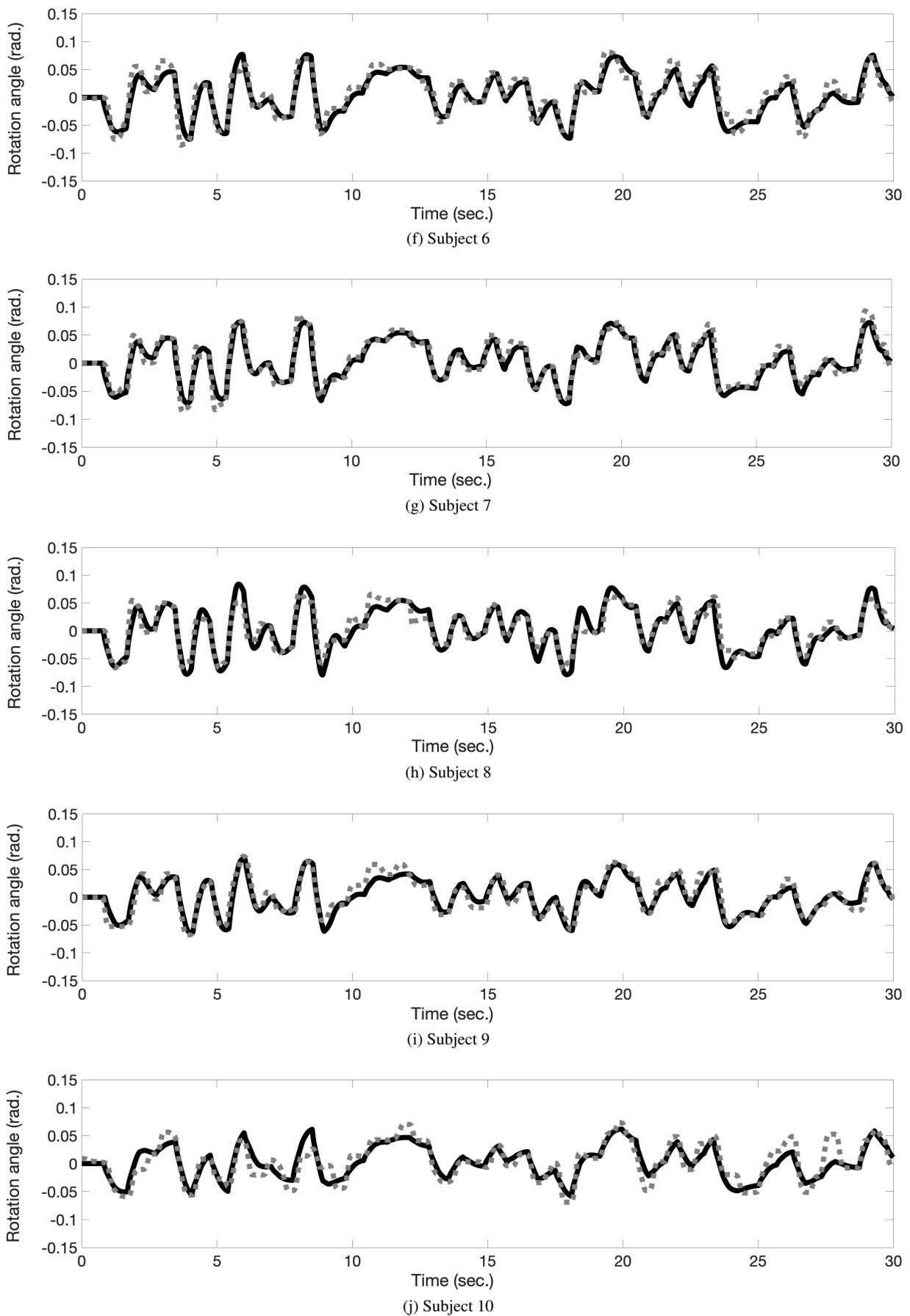


Fig. 4 The curve fitting with 10 experimental cases (continued). Solid lines represent estimated responses from the fitted models and dotted lines represent measured responses.

the L_1 -regularization is known to induce biased estimates, the bias with our method also slightly was decreased except for one parameter with the increased bias. The reason might be that our method pushes all parameters toward the mean of nominal parameter values obtained from 20 preliminary estimation. Therefore, if we set the appropriate values as the typical values, we could achieve lower values of bias and variance errors.

In the experimental study, our proposed method was compared with the Lasso in [20] using three performance criteria. First, we confirmed that the proposed method simultaneously selected and estimated sensitive parameters in the nonlinear model. As a result, five parameters were selected and estimated in our method, and the remaining parameters were fixed to the mean of nominal parameters obtained from preliminary estimation. With our method, 4 out of 5 sensitive parameters were the same as those selected in [20]. This result showed that our method behaved similarly in sensitive parameter selection by [20] using the Fisher information matrix. Selected sensitive parameters may vary slightly depending on the performance of the optimizer and the condition of initial points. However, they have not changed much over repeated randomized realizations. In addition, we presented VAF to quantitatively evaluate the goodness of fit of the estimated model. From the standard nonlinear least squares problem without L_1 -regularization, VAF was about 84.9%, and our method achieved about 82.5%. Hence, the goodness of fit of the model estimated from our method was similar to that of model estimated from the standard nonlinear least squares problem without L_1 -regularization although only 5 selected parameters were used for estimation in our method. Moreover, as shown in Fig. 4, most curve fitting errors occur around the peak points, which shows the limitation of the presented dynamics models that do not perfectly reflect the real physiological head-neck control processes. This could be due to the switching of human control strategies with sudden changes in head-neck orientations [2].

Next, the model identifiability was improved by sensitive parameter selection from our method. The variance of estimated 12 parameters is reduced by 71.1% on average. In general, when the number of parameters to be estimated is large with limited data, the model has poor (or lack of) identifiability. Therefore, through our method with key parameter selection, i.e., the nonlinear least squares problem with L_1 -regularization on the deviation of parameters from the mean of the nominal parameter values, the uniqueness of the estimated solution can be ensured even for an original problem with a lack of identifiability due to limited data.

Finally, the average computation time of our method was 54 times faster than that of the Lasso with the brute force optimization algorithm in [20]. Our method reduced the computation time by eliminating large inverse matrix computa-

tion by modifying a Jacobian formulation as a minimization formulation.

5 Conclusions

In this paper, we tackled a parameter estimation problem with limited data by formulating it as a nonlinear least squares estimator with L_1 -regularization.

We effectively improved the model identifiability by applying the Lasso to the nonlinear least squares problem. As asymptotic results, we provided consistency and oracle properties for a nonlinear regression approach with a generalized penalty function. Based on these results, we proposed a novel solution to our problem by solving the nonlinear least squares problem with L_1 -regularization on the deviation of parameters from the nominal values in order to simultaneously select and estimate model parameters.

From simulation and experimental studies, we successfully demonstrated that the proposed method simultaneously selected and estimated sensitive parameters, improved the model identifiability by reducing the variance of estimated parameters and took a much shorter computation time than that of the Lasso in [20].

Future work would be to apply our method to other clinical patient-specific calibration of the disease models [5,28].

Acknowledgements This work was supported by the Mid-career Research Programs through the National Research Foundation of Korea (NRF) funded by the Ministry of Science and ICT (NRF-2018R1A2B6008063, 2019R1A2C1002213). This publication was made possible by grant number U19AT006057 from the National Center for Complementary and Integrative Health (NCCIH) at the National Institutes of Health. Its contents are solely the responsibility of the authors and do not necessarily represent the official views of NCCIH. The research of Wei-Ying Wu was supported by Ministry of Science and Technology of Taiwan under grants (MOST 107-2118-M-259-001-).

Conflict of interest

The authors declare that they have no conflict of interest.

Appendix A Proof of Lemma 1

We first find the lower bound of $Q_n(\theta_0 + n^{-1/2}\mathbf{v}) - Q_n(\theta_0)$.

$$\begin{aligned} & Q_n(\theta_0 + n^{-1/2}\mathbf{v}) - Q_n(\theta_0) \\ &= S_n(\theta_0 + n^{-1/2}\mathbf{v}) - S_n(\theta_0) \\ & \quad + n\left(\sum_{k=1}^p p_{\lambda_n}(|\theta_{0k} + n^{-1/2}v_k|) - \sum_{k=1}^p p_{\lambda_n}(|\theta_{0k}|)\right) \\ &= n^{-1/2}(\nabla S_n(\theta_0))^T \mathbf{v} + \frac{1}{2}n^{-1}\mathbf{v}^T \nabla^2 S_n(\theta^*) \mathbf{v} \end{aligned}$$

$$\begin{aligned}
& + n \left(\sum_{k=1}^s p_{\lambda_n}(|\theta_{0k} + n^{-1/2}v_k|) - \sum_{k=1}^s p_{\lambda_n}(|\theta_{0k}|) \right) \\
& + n \sum_{k=s+1}^p p_{\lambda_n}(|\theta_{0k} + n^{-1/2}v_k|), \text{ where} \\
& \boldsymbol{\theta}^* = (\theta_1^*, \theta_2^*, \dots, \theta_p^*) \text{ lies between } \boldsymbol{\theta}_0 + n^{-1/2}\mathbf{v} \text{ and } \boldsymbol{\theta}_0, \\
& = n^{-1/2} (\nabla \mathbf{S}_n(\boldsymbol{\theta}_0))^T \mathbf{v} + \frac{1}{2} n^{-1} \mathbf{v}^T \nabla^2 \mathbf{S}_n(\boldsymbol{\theta}^*) \mathbf{v} \\
& + n \left(\sum_{k=1}^s q_{\lambda_n}(|\theta_{0k}^*|) \text{sgn}(\theta_{0k}^*) n^{-1/2} v_k \right) \\
& + n \sum_{k=s+1}^p p_{\lambda_n}(|\theta_k + n^{-1/2}v|) \\
& \geq n^{-1/2} (\nabla \mathbf{S}_n(\boldsymbol{\theta}_0))^T \mathbf{v} + \frac{1}{2} n^{-1} \mathbf{v}^T \nabla^2 \mathbf{S}_n(\boldsymbol{\theta}^*) (\mathbf{v}) \\
& + n \left(\sum_{k=1}^s q_{\lambda_n}(|\theta_{0k}^*|) \text{sgn}(\theta_{0k}^*) n^{-1/2} v_k \right) \\
& := \mathbb{A} + \mathbb{B} + \mathbb{C}.
\end{aligned}$$

For the sake of simplicity,

we define $\mathbf{d} = (d_1(\boldsymbol{\theta}, \boldsymbol{\theta}'), \dots, d_n(\boldsymbol{\theta}, \boldsymbol{\theta}'))$, where $d_i(\boldsymbol{\theta}, \boldsymbol{\theta}') = f(\mathbf{x}_i; \boldsymbol{\theta}) - f(\mathbf{x}_i; \boldsymbol{\theta}')$. For the term \mathbb{A} ,

$$\begin{aligned}
& n^{-1/2} (\nabla \mathbf{S}_n(\boldsymbol{\theta}_0))^T \mathbf{v} \\
& = -2n^{-1/2} \mathbf{v}^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \boldsymbol{\epsilon}, \\
& = -2\mathbf{v}^T \left[\mathbf{F}_n^{1/2} \left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2} \dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \boldsymbol{\epsilon} \right],
\end{aligned}$$

$$\text{where } \mathbf{F}_0 = \frac{1}{n} \dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0).$$

By Assumption 1, $\mathbf{F}_0^{1/2} \rightarrow \Gamma^{1/2}$ and we claim that,

$$\left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2} \dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \boldsymbol{\epsilon} \xrightarrow{d} N(0, \sigma^2 \mathbf{I}).$$

The claim follows from the lemma 2.1 in [8]. The condition for the lemma in our setting is

$$\left\| \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2} \right\|_{\infty} \rightarrow 0$$

where $\|A\|_{\infty}$ denote the maximum absolute value of all elements of matrix A . Since $\dot{\mathbf{F}}(\boldsymbol{\theta}_0)$ is an $n \times p$ matrix and $\left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2}$ is a $p \times p$ matrix,

$$\begin{aligned}
& \left\| \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2} \right\|_{\infty} \\
& \leq p \left\| \left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2} \right\|_{\infty} \left\| \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right\|_{\infty} \\
& = pn^{-1/2} \left\| \mathbf{F}_0^{-1/2} \right\|_{\infty} \left\| \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right\|_{\infty} \\
& = O\left(\frac{p}{\sqrt{n}}\right).
\end{aligned}$$

The first and second conditions from Assumption 1 imply the last equality. Therefore, we have

$$\begin{aligned}
& -2\mathbf{v}^T \left[\mathbf{F}_0^{1/2} \left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)^{-1/2} \dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \boldsymbol{\epsilon} \right] \\
& \xrightarrow{d} N(0, 4\sigma^2 \mathbf{v}^T \Gamma \mathbf{v}).
\end{aligned} \tag{1}$$

For the term \mathbb{B} ,

$$\begin{aligned}
& \frac{1}{2} n^{-1} \mathbf{v}^T \nabla^2 \mathbf{S}_n(\boldsymbol{\theta}^*) \mathbf{v} \\
& = n^{-1} \mathbf{v}^T \left\{ \left(\dot{\mathbf{F}}(\boldsymbol{\theta}^*)^T \dot{\mathbf{F}}(\boldsymbol{\theta}^*) \right) + \ddot{\mathbf{F}}(\boldsymbol{\theta}^*)^T (\mathbf{I} \otimes (\mathbf{d} - \boldsymbol{\epsilon})) \right\} \mathbf{v} \\
& = \mathbf{v}^T \left[\frac{\left(\dot{\mathbf{F}}(\boldsymbol{\theta}_0)^T \dot{\mathbf{F}}(\boldsymbol{\theta}_0) \right)}{n} \mathbf{Z}_n^{-1} \mathbf{v} \right], \text{ where} \\
& \mathbf{Z}_n = \left\{ \left(\dot{\mathbf{F}}(\boldsymbol{\theta}^*)^T \dot{\mathbf{F}}(\boldsymbol{\theta}^*) \right) + \ddot{\mathbf{F}}(\boldsymbol{\theta}^*)^T (\mathbf{I} \otimes (\mathbf{d} - \boldsymbol{\epsilon})) \right\}^{-1} n \mathbf{F}_0.
\end{aligned}$$

If we show $\mathbf{Z}_n^{-1} \xrightarrow{p} \mathbf{I}$, with Assumption 1, we obtain

$$\mathbb{A} + \mathbb{B} \xrightarrow{d} N(\mathbf{v}^T \Gamma \mathbf{v}, 4\sigma^2 \mathbf{v}^T \Gamma \mathbf{v}). \tag{2}$$

\mathbf{Z}_n^{-1} can be rewritten as,

$$\begin{aligned}
\mathbf{Z}_n^{-1} & = (n \mathbf{F}_0)^{-1} \left\{ \left(\dot{\mathbf{F}}(\boldsymbol{\theta}^*)^T \dot{\mathbf{F}}(\boldsymbol{\theta}^*) \right) + \ddot{\mathbf{F}}(\boldsymbol{\theta}^*)^T (\mathbf{I} \otimes (\mathbf{d} - \boldsymbol{\epsilon})) \right\} \\
& = (n \mathbf{F}_0)^{-1} \left(\dot{\mathbf{F}}(\boldsymbol{\theta}^*)^T \dot{\mathbf{F}}(\boldsymbol{\theta}^*) \right) + (n \mathbf{F}_0)^{-1} \ddot{\mathbf{F}}(\boldsymbol{\theta}^*)^T (\mathbf{I} \otimes \mathbf{d}) \\
& \quad - (n \mathbf{F}_0)^{-1} \ddot{\mathbf{F}}(\boldsymbol{\theta}^*)^T (\mathbf{I} \otimes \boldsymbol{\epsilon}).
\end{aligned}$$

By Assumption 1, the first term converges to \mathbf{I} almost surely. By the conditions 1, 2, and 4 of Assumption 1 and Cauchy-Schwarz inequality, the second term converges to zero almost surely. For the last term, it is enough to show that

$$\frac{1}{n} \mathbf{f}_{ks}(\boldsymbol{\theta})^T \boldsymbol{\epsilon} \xrightarrow{p} 0 \tag{3}$$

uniformly on $S - \{\theta \in \Theta : \|\theta - \theta_0\| \leq \delta\}$ with probability 1, because of the conditions 1 and 2 of Assumption 1. Then, by the condition 4 of Assumption 1, (3.17) can be shown in a manner similar to that of [26] For the term \mathbb{C} , by Assumption 2 for a fixed $s < \infty$,

$$\left(n^{1/2} \sum_{k=1}^s q_{\lambda_n}(|\theta_{0k}^*|) \text{sgn}(\theta_{0k}^*) v_k \right) \rightarrow 0. \tag{4}$$

Thus, combined (2) with (4), we have

$$\mathbb{A} + \mathbb{B} + \mathbb{C} \xrightarrow{d} N(\mathbf{v}^T \Gamma \mathbf{v}, 4\sigma^2 \mathbf{v}^T \Gamma \mathbf{v}),$$

which leads to the desired result for large enough C .

Appendix B Proof of Theorem 1

By Lemma 1 and the continuity of $\mathbf{Q}_n(\cdot)$, we obtain Theorem 1.

Appendix C Proof of Theorem 2

Proof of (i)

First, we break M_k into two sets:

$$\begin{aligned} M_k &= \left\{ \omega : \hat{\theta}_k \neq 0, |\hat{\theta}_k| \geq Cn^{-1/2} \right\} \\ &\quad + \left\{ \omega : \hat{\theta}_k \neq 0, |\hat{\theta}_k| < Cn^{-1/2} \right\} \\ &=: E_n + F_n. \end{aligned}$$

Then, it is enough to show for any $\epsilon > 0$, $P(E_n) < \epsilon/2$ and $P(F_n) < \epsilon/2$. For any $\epsilon > 0$, we can show $P(E_n) < \epsilon/2$ for large enough n because of the consistency.

To verify $P(F_n) < \epsilon/2$ for large enough n , we first show $n^{1/2}q_{\lambda_n}(|\hat{\theta}_k|) = O_p(1)$ on the set F_n . Note that

$$\begin{aligned} &n^{-1/2}\nabla S_n(\boldsymbol{\theta}) - n^{-1/2}\nabla S_n(\boldsymbol{\theta}_0) \\ &= n^{-1/2}\nabla^2 S_n(\boldsymbol{\theta}^{**})(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\ &= n^{1/2} \frac{(\dot{\mathbf{F}}(\boldsymbol{\theta}^{**})^T \dot{\mathbf{F}}(\boldsymbol{\theta}^{**}))}{n} \mathbf{Z}_n^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = O_p(1), \end{aligned}$$

where $\boldsymbol{\theta}^{**}$ lies between $\boldsymbol{\theta}$ and $\boldsymbol{\theta}_0$. Since $\frac{1}{n}\dot{\mathbf{F}}(\boldsymbol{\theta}^{**})^T \dot{\mathbf{F}}(\boldsymbol{\theta}^{**}) \xrightarrow{p} \Gamma$, $\mathbf{Z}_n^{-1} \xrightarrow{p} I$ and $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| = O_p(n^{-1/2})$. Thus, we have

$$\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \leq Cn^{-1/2}} \|n^{-1/2}\nabla S_n(\boldsymbol{\theta}) - n^{-1/2}\nabla S_n(\boldsymbol{\theta}_0)\| = O_p(1). \quad (1)$$

Combining (1) with $\|n^{-1/2}\nabla S_n(\boldsymbol{\theta}_0)\| = O_p(1)$, we have

$$\|n^{-1/2}\nabla S_n(\boldsymbol{\theta})\| = O_p(1) \quad (2)$$

for $\boldsymbol{\theta}$ which satisfies $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \leq Cn^{-1/2}$. Since $\hat{\boldsymbol{\theta}}$ is the local minimizer of $\mathbf{Q}_n(\boldsymbol{\theta})$ with the root-n-consistent, we attain

$$n^{1/2}q_{\lambda_n}(|\hat{\theta}_k|) = O_p(1) \quad (3)$$

from

$$\begin{aligned} &n^{-1/2} \frac{\partial \mathbf{Q}_n(\boldsymbol{\theta})}{\partial \theta_k} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ &= n^{-1/2} \frac{\partial S_n(\boldsymbol{\theta})}{\partial \theta_k} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} + n^{1/2}q_{\lambda_n}(|\hat{\theta}_k|)\text{sgn}(\hat{\theta}_k). \end{aligned}$$

Therefore, there exists M' such that, for large enough n ,

$$P\{\omega : \hat{\theta}_k \neq 0, |\hat{\theta}_k| < Cn^{-1/2}, n^{1/2}q_{\lambda_n}(|\hat{\theta}_k|) > M'\} < \epsilon/2.$$

In addition, by the second assumption of Assumption 2,

$$\begin{aligned} &\{\omega : \hat{\theta}_k \neq 0, |\hat{\theta}_k| < Cn^{-1/2}, n^{1/2}q_{\lambda_n}(|\hat{\theta}_k|) > M'\} \\ &= \{\omega : \hat{\theta}_k \neq 0, |\hat{\theta}_k| < Cn^{-1/2}\} \end{aligned}$$

for the large enough n , which leads to $P(F_n) < \epsilon/2$ for large enough n .

Proof of (ii)

By the Taylor expansion,

$$\begin{aligned} n^{-1/2}\nabla \mathbf{Q}_n(\hat{\boldsymbol{\theta}}) &= n^{-1/2}\nabla S_n(\hat{\boldsymbol{\theta}}) + n^{-1/2}\mathbf{q}_{\lambda_n}(\hat{\boldsymbol{\theta}}) \cdot \text{sgn}(\hat{\boldsymbol{\theta}}) \\ &= n^{-1/2} \left(\nabla S_n(\boldsymbol{\theta}_0) + \nabla^2 S_n(\boldsymbol{\theta}^{**})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \right) \\ &\quad + n^{1/2}\mathbf{q}_{\lambda_n}(\hat{\boldsymbol{\theta}}) \cdot \text{sgn}(\hat{\boldsymbol{\theta}}) \end{aligned}$$

where $\mathbf{q}_{\lambda_n}(\hat{\boldsymbol{\theta}}) \cdot \text{sgn}(\hat{\boldsymbol{\theta}}) = \left(q_{\lambda_n}(|\hat{\theta}_1|)\text{sgn}(\hat{\theta}_1), \dots, q_{\lambda_n}(|\hat{\theta}_p|)\text{sgn}(\hat{\theta}_p) \right)^T$. Since $\hat{\boldsymbol{\theta}}$ is the local minimizer of $\mathbf{Q}_n(\boldsymbol{\theta})$, $\nabla \mathbf{Q}_n(\hat{\boldsymbol{\theta}}) = 0$ so that

$$\begin{aligned} n^{-1/2}(-\nabla S_n(\boldsymbol{\theta}_0)) &= n^{-1}\nabla^2 S_n(\boldsymbol{\theta}^{**}) \left(n^{1/2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \right) \\ &\quad + n^{1/2}\mathbf{q}_{\lambda_n}(\hat{\boldsymbol{\theta}}) \cdot \text{sgn}(\hat{\boldsymbol{\theta}}). \end{aligned}$$

Finally,

$$n^{1/2}2\Gamma_{11}(\hat{\boldsymbol{\theta}}_{11} - \boldsymbol{\theta}_{01} + (2\Gamma_{11})^{-1}b_n) \xrightarrow{d} N(0, 4\Gamma_{11}\sigma^2),$$

because $n^{-1}\nabla^2 S_n(\boldsymbol{\theta}^{**}) \xrightarrow{p} 2\Gamma$, $n^{-1/2}(-\nabla S_n(\boldsymbol{\theta}_0)) \xrightarrow{d} N(\mathbf{0}, 4\Gamma\sigma^2)$ and the consistency of $\hat{\boldsymbol{\theta}}$.

Appendix D Nonlinear least squares estimator with L_1 -regularization.

In order to apply the Lasso to the nonlinear least squares problem, we reformulate the Levenberg-Marquardt (LM) optimization algorithm as a linear least squares problem as follow.

$$\begin{aligned} \hat{\boldsymbol{\theta}}^{\hat{j}+1} &= \min_{\boldsymbol{\theta}} \|\boldsymbol{\Lambda}\tilde{\boldsymbol{\theta}}^{\hat{j}+1} - \boldsymbol{\Lambda}\tilde{\boldsymbol{\theta}}^{\hat{j}} + \mathbf{J}^T \tilde{\mathbf{r}}(\tilde{\boldsymbol{\theta}}^{\hat{j}})\|_2^2 \\ &\text{s.t. } \sum_{k=1}^p |\tilde{\theta}_k^{\hat{j}+1}| \leq T_{\theta}, \end{aligned} \quad (1)$$

where $\tilde{\boldsymbol{\theta}}^{\hat{j}+1} = \tilde{\boldsymbol{\theta}}^{\hat{j}} + \Delta\boldsymbol{\theta}$, $\Delta\boldsymbol{\theta} = -\boldsymbol{\Lambda}^{-1}\mathbf{J}^T \tilde{\mathbf{r}}(\tilde{\boldsymbol{\theta}}^{\hat{j}})$, $\boldsymbol{\Lambda} = (\mathbf{J}^T \mathbf{J} + \mu \text{diag}(\mathbf{J}^T \mathbf{J}))$. The residual vector $\tilde{\mathbf{r}}(\tilde{\boldsymbol{\theta}}^{\hat{j}}) = \mathbf{r}(\tilde{\boldsymbol{\theta}}^{\hat{j}} + \boldsymbol{\theta}) = \mathbf{r}(\boldsymbol{\theta}^{\hat{j}}) = \mathbf{y} - f(\mathbf{X}; \boldsymbol{\theta}^{\hat{j}})$. T_{θ} is the regularization hyperparameter. The sum of the deviation of parameters from the mean of the nominal parameter values is less than or equal to the regularization hyperparameter T_{θ} . The damping factor μ affects the efficiency and the convergence stability [3]. \mathbf{J} is the Jacobian matrix which consists of all first-order partial derivatives of the residual vector with respect to the parameters, evaluated for $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}^{\hat{j}}$ as follows.

$$\mathbf{J} := \frac{\partial \tilde{\mathbf{r}}(\tilde{\boldsymbol{\theta}}^{\hat{j}})}{\partial \tilde{\boldsymbol{\theta}}^{\hat{j}}}$$

In addition, as compared with the conventional Lasso fixing insensitive parameters to 0, our method simultaneously selects and estimates only sensitive parameters while fixing insensitive parameters onto the mean of the nominal parameter values $\bar{\boldsymbol{\theta}}$.

Table 3 The algorithm of sensitive parameter selection using our method

Input:	(1) Experimental Data and Dynamical Model (2) Vector of Normalized Values ($\tilde{\Phi}$), where the Mean of Nominal Parameter Values is $\bar{\theta}$ (3) Desired Optimal Number of Sensitive Parameters (n^*) (4) The Regularization Hyperparameter T_θ
Output:	(1) A Subset of the Sensitive Parameters ($\hat{\theta}$)
<pre> 1: NumParams = 0 2: $T_\theta = 1.0$ 3: while NumParams $\neq n^*$ do 4: repeat 5: Solve (1) in Appendix D 6: until $\tilde{\Phi}$ convergences 7: for $i = 1 : n$ do 8: if $\hat{\Phi}(i) > 0.001$ then 9: NumParams = NumParams + 1 10: end if 11: end for 12: if NumParams $\neq n^*$ then 13: $T_\theta = T_\theta - \frac{\text{NumParams} - n^*}{i + n^*}$ 14: end if 15: end while 16: $\hat{\theta}$, where a subset of the normalized sensitive parameters is $\hat{\Phi}$ </pre>	

References

- Babyak, M.A.: What you see may not be what you get: a brief, nontechnical introduction to overfitting in regression-type models. *Psychosomatic medicine* **66**(3), 411–421 (2004)
- Chen, K.J., Keshner, E., Peterson, B., Hain, T.: Modeling head tracking of visual targets. *Journal of Vestibular Research* **12**(1), 25–33 (2002)
- Cui, M., Zhao, Y., Xu, B., Gao, X.w.: A new approach for determining damping factors in levenberg-marquardt algorithm for solving an inverse heat conduction problem. *International Journal of Heat and Mass Transfer* **107**, 747–754 (2017)
- Do, H.N., Choi, J., Lim, C.Y., Maiti, T.: Appearance-based localization of mobile robots using group lasso regression. *Journal of Dynamic Systems, Measurement, and Control* **140**(9), 091016 (2018)
- Do, H.N., Ijaz, A., Gharahi, H., Zambrano, B., Choi, J., Lee, W., Baek, S.: Prediction of abdominal aortic aneurysm growth using dynamical gaussian process implicit surface. *IEEE Transactions on Biomedical Engineering* (2018)
- Fan, J., Li, R.: Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* **96**(456), 1348–1360 (2001)
- Gutenkunst, R.N., Waterfall, J.J., Casey, F.P., Brown, K.S., Myers, C.R., Sethna, J.P.: Universally sloppy parameter sensitivities in systems biology models. *PLoS computational biology* **3**(10), e189 (2007)
- Huber, P.J., et al.: Robust regression: asymptotics, conjectures and monte carlo. *The Annals of Statistics* **1**(5), 799–821 (1973)
- Johnson, B.A., Lin, D., Zeng, D.: Penalized estimating functions and variable selection in semiparametric regression models. *Journal of the American Statistical Association* **103**(482), 672–680 (2008)
- Jolliffe, I.T., Trendafilov, N.T., Uddin, M.: A modified principal component technique based on the lasso. *Journal of computational and Graphical Statistics* **12**(3), 531–547 (2003)
- Kump, P., Bai, E.W., Chan, K.S., Eichinger, B., Li, K.: Variable selection via rival (removing irrelevant variables amidst lasso iterations) and its application to nuclear material detection. *Automatica* **48**(9), 2107–2115 (2012)
- Kyubaek, Y., Jongeun, C.: Penalized nonlinear regression with application to head-neck position tracking. In: *ASME 2019 Dynamic Systems and Control Conference*. American Society of Mechanical Engineers Digital Collection (in press)
- Little, M.P., Heidenreich, W.F., Li, G.: Parameter identifiability and redundancy: theoretical considerations. *PLoS one* **5**(1), e8915 (2010)
- Lund, A., Dyke, S.J., Song, W., Billionis, I.: Global sensitivity analysis for the design of nonlinear identification experiments. *Nonlinear Dynamics* **98**(1), 375–394 (2019)
- Moré, J.J.: The levenberg-marquardt algorithm: implementation and theory. In: *Numerical analysis*, pp. 105–116. Springer (1978)
- Peng, G., Hain, T., Peterson, B.: A dynamical model for reflex activated head movements in the horizontal plane. *Biological cybernetics* **75**(4), 309–319 (1996)
- Pitt, M.A., Myung, I.J.: When a good fit can be bad. *Trends in cognitive sciences* **6**(10), 421–425 (2002)
- Popovich Jr, J.M., Reeves, N.P., Priess, M.C., Cholewicki, J., Choi, J., Radcliffe, C.J.: Quantitative measures of sagittal plane head-neck control: A test-retest reliability study. *Journal of biomechanics* **48**(3), 549–554 (2015)
- Qin, P., Nishii, R., Yang, Z.J.: Selection of narx models estimated using weighted least squares method via gic-based method and l1-norm regularization methods. *Nonlinear Dynamics* **70**(3), 1831–1846 (2012)
- Ramadan, A., Boss, C., Choi, J., Reeves, N.P., Cholewicki, J., Popovich, J.M., Radcliffe, C.J.: Selecting sensitive parameter subsets in dynamical models with application to biomechanical system identification. *Journal of biomechanical engineering* **140**(7), 074503 (2018)
- Ramadan, A., Choi, J., Cholewicki, J., Reeves, N.P., Popovich, J.M., Radcliffe, C.J.: Feasibility of incorporating test-retest reliability and model diversity in identification of key neuromuscular pathways during head position tracking. *IEEE transactions on neural systems and rehabilitation engineering* (2019)
- Rasouli, M., Westwick, D., Rosehart, W.: Reducing induction motor identified parameters using a nonlinear lasso method. *Electric Power Systems Research* **88**, 1–8 (2012)
- Tibshirani, R.: Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 267–288 (1996)
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., Knight, K.: Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **67**(1), 91–108 (2005)
- Tibshirani, R., Wasserman, L.A.: Sensitive parameters. *Canadian Journal of Statistics* **16**(2), 185–192 (1988)
- Wu, C.F.: Asymptotic theory of nonlinear least squares estimation. *The Annals of Statistics* pp. 501–513 (1981)
- Yuan, M., Lin, Y.: Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**(1), 49–67 (2006)
- Zhang, L., Jiang, Z., Choi, J., Lim, C.Y., Maiti, T., Baek, S.: Patient-specific prediction of abdominal aortic aneurysm expansion using bayesian calibration. *IEEE Journal of Biomedical and Health Informatics* (2019)
- Zou, H.: The adaptive lasso and its oracle properties. *Journal of the American statistical association* **101**(476), 1418–1429 (2006)

-
30. Zou, H., Hastie, T.: Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **67**(2), 301–320 (2005)