

# Group LASSO Approach for Estimation and Source Localization of Forced Oscillations

Rajasekhar Anguluri, *Member, IEEE*, Lalitha Sankar, *Senior Member, IEEE*, and Oliver Kosut, *Member, IEEE*

**Abstract**—We consider a group LASSO (linear absolute shrinkage and selection operator) based estimator for jointly estimating the initial state and multiple forced oscillations (FOs), including their locations, in a power system. Through rigorous theoretical arguments, our main goal in this paper is to convey that the proposed estimator is desirable when the number of sources of FOs is sparse. Based on the observability and forced response matrices of a linearized power system model, we first develop a mutual incoherence condition. Using this condition, we derive  $\ell_2$ -error bound on the estimation accuracy and group support recovery condition for accurately locating the sources of FO. We also establish relationship between the group LASSO based location recovery condition to that of left invertibility condition—a classical concept that determines the existence of unknown input estimators for linear systems. Finally, we validate the performance of the proposed approach via numerical simulations on IEEE 68-bus and WECC 179-bus power systems.

**Index Terms**—Forced oscillations, group LASSO, source location, estimation.

## I. INTRODUCTION

### Related work:

- Dissipating Energy Method (Chavlier, Senroy, Mani)
- Data-Driven approaches (Huang, etc)
- Statistical approaches
- frequency based approaches

The main contributions of this paper are as follows:

- 1) Assuming a linearized electro-mechanical model for a power system, we propose a group LASSO optimization program to simultaneously estimate the initial state and multiple FO inputs, and to recover the locations of FOs.
- 2) Based on the mutual incoherence condition for the joint observability and forced response matrix, we provide support recovery condition for group LASSO to accurately locate the sources of FO with high probability, and derive probabilistic bound on the estimation error.
- 3) We establish connection between group LASSO based support recovery conditions to that of the left-invertibility condition—a classical concept in dynamical system inversion literature—which guarantees the existence of estimators for recovering unknown inputs in linear systems.
- 4) Using the proposed group LASSO approach, we demonstrate that multiple forced oscillations present in IEEE 68-bus and WECC 179-bus benchmark power systems can be efficiently estimated including their source locations.

The remainder of the paper is organized as follows. Section II introduces the power system model and description of the group lasso framework. In Section III and IV, we state and prove consistency results on estimation and source location of forced oscillations, respectively. In Section V, we present the results of simulations that validate our theoretical analysis. Section VI concludes the paper.

## II. PROBLEM SETUP AND PRELIMINARY NOTIONS

### A. Power System Dynamics and Measurement Models

For a multi-machine power system, the small-signal model near an operating point can be modelled using the following continuous time invariant linear dynamical system:

$$\dot{\mathbf{x}}(t) = \tilde{\mathbf{A}}\mathbf{x}(t) + \tilde{\mathbf{B}}\tilde{\mathbf{u}}(t), \quad (1)$$

where  $\mathbf{x}(t) \in \mathbb{R}^n$  is the state vector (generator rotor angle and speed deviations) and  $\tilde{\mathbf{A}}$  is the state matrix of appropriate dimension. The forced oscillations input  $\tilde{\mathbf{u}}(t) = [\tilde{u}_1(t), \dots, \tilde{u}_m(t)] \in \mathbb{R}^m$  is injected by  $m$  potentially malfunctioned<sup>1</sup> generators, whose locations are encoded in the matrix  $\tilde{\mathbf{B}}$ . Although in the FOs source localization literature  $\mathbf{u}(t)$  is assumed to be periodic, we make no assumptions about its periodicity, except for being a deterministic signal.

Using the standard state-space discretization technique [], the discrete system dynamics of (1) along with the measurement model is given by

$$\begin{aligned} \mathbf{x}[k+1] &= \mathbf{A}\mathbf{x}[k] + \mathbf{B}\tilde{\mathbf{u}}[k], \\ \mathbf{y}[k] &= \mathbf{C}\mathbf{x}[k] + \mathbf{v}[k], \end{aligned} \quad (2)$$

where  $\mathbf{y}[k] \in \mathbb{R}^r$  is the noisy measurement of the state  $\mathbf{x}[k]$ , with noise  $\mathbf{v}[k] \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . Further,  $\mathbf{A} = e^{\tilde{\mathbf{A}}\Delta t}$  and  $\mathbf{B} = (\int_{\tau=0}^{\Delta t} e^{\tilde{\mathbf{A}}t} dt)\tilde{\mathbf{B}}$ , where  $\Delta t$  is the sampling time, and the output matrix  $\mathbf{C} = [\mathbf{e}_1, \dots, \mathbf{e}_r]^T$  describes PMU locations in the network. Finally, we assume that the initial state  $\mathbf{x}[0] := \mathbf{x}_0$  is unknown.

Let  $\mathbf{y} = [\mathbf{y}^T[1] \dots \mathbf{y}^T[N]]^T \in \mathbb{R}^{Nr}$  denote the measurements collected over the time horizon  $k = 1, \dots, N$ . Analogously, define  $\tilde{\mathbf{u}} = [\tilde{\mathbf{u}}^T[0] \dots \tilde{\mathbf{u}}^T[N-1]]^T \in \mathbb{R}^{Nm}$  and  $\mathbf{v} = [\mathbf{v}^T[1] \dots \mathbf{v}^T[N]]^T \in \mathbb{R}^{Nr}$ . Then, from (2)  $\mathbf{y}$  can be expressed as

$$\mathbf{y} = \mathbf{O}\mathbf{x}_0 + \mathbf{J}\tilde{\mathbf{u}} + \mathbf{v}, \quad (3)$$

The authors are with the School of Electrical, Computer, and Energy Engineering, Arizona State University, Tempe, AZ 85281 USA (e-mail: {rangulur,lalithasankar,okosut}@asu.edu).

<sup>1</sup>Here, the term "malfunction" refers to a specific failed component of the generator (e.g., governor control, excitation system etc.), but not the overall generator.

where the observability matrix  $\mathbf{O} \in \mathbb{R}^{Nr \times n}$  and the impulse response matrix  $\tilde{\mathbf{J}} \in \mathbb{R}^{Nr \times Nm}$  are defined as

$$\mathbf{O} := \begin{bmatrix} \mathbf{CA} \\ \mathbf{CA}^2 \\ \vdots \\ \mathbf{CA}^N \end{bmatrix}; \tilde{\mathbf{J}} := \begin{bmatrix} \mathbf{CB} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{CAB} & \mathbf{CB} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{CA}^{N-1}\mathbf{B} & \mathbf{CA}^{N-2}\mathbf{B} & \dots & \mathbf{CB} \end{bmatrix}.$$

### B. Group LASSO for Initial State and Input Estimation

We assume that  $m_0$  out of  $m$  sources are injecting the forced inputs in to the system (2), where  $m_0 \ll m$ . This is a valid assumption, since in many widely reported incidents on FOs, only a few malfunctioned controllers (one for every generator) turned out to be the sources.

In order to capture this sparsity structure among the source locations, we reshape the input  $\tilde{\mathbf{u}}$  as a concatenation of  $m$  subvectors, where each subvector contain the time stacked inputs at a given location. Formally, let  $\mathbf{u}_j = [\tilde{u}_j[0], \dots, \tilde{u}_j[N-1]]^T$  denote the input injected by the  $j$ -th source,  $j \in \{1, \dots, m\}$ , over a time interval of length  $N$ . Then, one can define a permutation/reshaping matrix  $\mathbf{P}$  such that  $\tilde{\mathbf{u}} = \mathbf{P}[\mathbf{u}_1^T, \dots, \mathbf{u}_m^T]^T$ . Define  $\mathbf{J} := \tilde{\mathbf{J}}\mathbf{P} = [\mathbf{J}_1, \dots, \mathbf{J}_m]$ . To emphasize that the vector  $(\mathbf{x}_0, \mathbf{u})$  is the true parameter, we denote it by  $(\mathbf{x}_0^*, \mathbf{u}^*)$ , where  $\mathbf{u}^* = [(\mathbf{u}_1^*)^T, \dots, (\mathbf{u}_m^*)^T]^T$ . Then, from (3) we have

$$\begin{aligned} \mathbf{y} &= \mathbf{O}\mathbf{x}_0^* + \tilde{\mathbf{J}}\mathbf{P}\mathbf{u}^* + \mathbf{v} \\ &= \mathbf{O}\mathbf{x}_0^* + \sum_{j=1}^m \mathbf{J}_j \mathbf{u}_j^* + \mathbf{v}. \end{aligned} \quad (4)$$

Since  $\mathbf{u}_j^* \neq \mathbf{0}$  if and only if the  $j$ -th source injects input, from (4) and the sparsity assumption ( $m_0 \ll m$ ), it is clear that only few non-zero  $\beta_j$ 's contribute to the measurement  $\mathbf{y}$ . Finally, notice that  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ .

We now state the group LASSO optimization problem to estimate the group sparse vector  $(\mathbf{x}_0^*, \mathbf{u}^*)$ :

$$\begin{bmatrix} \hat{\mathbf{x}}_0 \\ \hat{\mathbf{u}} \end{bmatrix} = \arg \min_{\substack{\mathbf{x}_0 \in \mathbb{R}^n \\ \mathbf{u} \in \mathbb{R}^{Nm}}} \left\{ \frac{1}{2T} \|\mathbf{y} - \mathbf{O}\mathbf{x}_0 - \mathbf{J}\mathbf{u}\|_2^2 + \lambda_T \sum_{j=1}^m \|\mathbf{u}_j\|_2 \right\}. \quad (5)$$

Here,  $T = Nm$  denotes the dimension of  $\mathbf{y}$ , and the tuning parameter  $\lambda_T > 0$  penalizes unknown parameters in the model (4). We make few comments about the optimization program (5). First, it is a convex optimization problem, and there exist several computationally efficient algorithms for obtaining  $\hat{\beta}$ , including the coordinated LARS algorithm. We provide more details about numerical algorithms in Section V. Second, although the penalty is over the input vector  $\mathbf{u}$ , the optimization framework jointly recovers the initial state and the input. Third, due to the mixed  $\ell_1/\ell_2$ -norm penalty, i.e., sum of  $\ell_2$  norm of the subvectors  $\mathbf{u}_j$ s, the estimates  $\hat{\mathbf{u}}_j$  including  $\hat{\mathbf{x}}_0$  are not sparse; see [1]. Thus, the estimates returned by (5) are useful in applications where the vector parameters are fewer in number but each vector is dense—e.g., initial state deviations and sustained oscillations in power systems

★ **WRITE ABOUT IDENTIFIABILITY ASSUMPTION/CONDITION**  
**ADD FIGURE/BLOCK DIAGRAM OF THE ESTIMATION SETUP**

★

We quantify the performance of an estimator  $\hat{\beta} = (\hat{\mathbf{x}}_0, \hat{\mathbf{u}})$  using the following error metrics:

- $\hat{\beta}$  is said to be  $\ell_2$ -consistent if  $\|\hat{\beta} - \beta^*\|_2 \leq o(T)$  with probability at least  $1 - c_1 \exp(-c_2 T)$ , for some constants  $c_1, c_2 > 0$ .
- $\hat{\mathbf{u}}$  is said to be group selection consistent if  $\text{gsupp}(\hat{\mathbf{u}}) = \text{gsupp}(\mathbf{u}^*)$  with probability at least  $1 - c_1 \exp(-c_2 T)$ , for some constants  $c_1, c_2 > 0$ ; here,  $\text{gsupp}(\mathbf{u}) = \{j : \mathbf{u}_j \neq \mathbf{0}, j = 1, \dots, m\}$ .

Here, the  $\ell_2$ -error bound ensures that the joint estimate  $(\hat{\mathbf{x}}_0, \hat{\mathbf{u}})$  can be made arbitrarily close to the true parameter  $(\mathbf{x}_0^*, \mathbf{u}^*)$  by increasing the number of samples  $T = Nm$ . Instead, the group selection consistency guarantees that as  $T \rightarrow \infty$ ,  $\hat{\mathbf{u}}$  correctly identifies the sources of FOs. **non-asymptotic nature** In the following sections we derive sufficient conditions, which guarantee that  $\hat{\beta}$  obtained from (5) is both  $\ell_2$ -estimation and group selection consistent. Moreover, our results also provide lower bounds on the measurements  $T$  required to achieve certain level of performance.

### III. INITIAL STATE AND INPUT RECOVERY: ESTIMATION CONSISTENCY

Let  $\tilde{\mathcal{G}} = \{1, \dots, m\}$  be the indices of FO source locations, and  $\mathcal{G} = \tilde{\mathcal{G}} \cup \{0\}$ .

### IV. INITIAL STATE AND INPUT RECOVERY: LOCATION SELECTION CONSISTENCY

**Assumption IV.1. (Group normalization and identifiability conditions)** Let  $\mathbf{M}_S = [\mathbf{O} \quad \mathbf{J}_S]$ , where  $\mathbf{J}_S$  is a sub-matrix of  $\mathbf{J} = [\mathbf{J}_1, \dots, \mathbf{J}_m]$  constructed by concatenating column blocks  $\mathbf{J}_j$ ,  $j \in S$ .

(A1) The column block matrices  $\mathbf{O}$  and  $\mathbf{J}_j$  satisfies the group normalization condition, i.e.,

$$\max \left\{ \frac{\|\mathbf{O}\|_2}{\sqrt{T}}, \frac{\|\mathbf{J}_1\|_2}{\sqrt{T}}, \dots, \frac{\|\mathbf{J}_m\|_2}{\sqrt{T}} \right\} \leq C. \quad (6)$$

(A2) The minimum singular value of the scaled submatrix  $\mathbf{M}_S$  is bounded below:

$$\gamma_{\min} \left( \frac{1}{T} \mathbf{M}_S^T \mathbf{M}_S \right) \geq c_{\min} > 0. \quad (7)$$

(A3) There exists some  $\alpha \in (0, 1)$ , referred to as "block-mutual incoherence" parameter, such that

$$\max_{j \in S^c} \left\| (\mathbf{J}_j^T \mathbf{M}_S) (\mathbf{M}_S^T \mathbf{M}_S)^{-1} \right\|_2 \leq 1 - \alpha. \quad (8)$$

Let  $\Gamma_j$  satisfies  $\mathbf{u}_j^* = \Gamma_j [(\mathbf{x}_0^*)^T \quad (\mathbf{u}^*)^T]^T$ , for any  $j \in S$ . Define

$$f(\mathbf{M}_S) = \max_{j \in S} \sum_{l \in S} \left\| \Gamma_j (\mathbf{M}_S^T \mathbf{M}_S)^{-1} \Gamma_l^T \right\|_2 \quad (9)$$

The above function extend matrix  $\|\cdot\|_\infty$  norm to the case of block matrix; here  $\|A\|_\infty = \max_{i=1, \dots, s} |A_{i,j}|$ .

**Theorem IV.2. (Location recovery)** Consider the linear model (4) satisfying assumptions (A1)–(A3). Suppose that the regularization parameter  $\lambda_T$  satisfies

$$\lambda_T \geq \frac{2C\sigma}{(1-\alpha)} \left( \sqrt{\frac{N}{T}} + \sqrt{\frac{\log(m-m_0)}{T}} \right). \quad (10)$$

Then with probability at least  $1 - \frac{2m}{m_0(m-m_0)}$ , the group lasso estimate (5) has the following properties:

- 1) (No false inclusion) The optimal solution  $\hat{\mathbf{u}}$  is unique with its group support contained within the true support  $S$  (i.e.,  $S(\hat{\mathbf{u}}) \subseteq S(\mathbf{u}^*)$ ).
- 2) (No false exclusion) Let

$$g(\lambda_T) = \frac{2\sigma}{\sqrt{c_{\min}}} \left( \sqrt{\frac{N}{T}} + \sqrt{\frac{\log m_0}{T}} \right) + \tilde{f}(\mathbf{M}_S) \lambda_T T, \quad (11)$$

If the minimum  $\ell_2$  norm of the true input  $\mathbf{u}^*$  on its group support  $S$  is bounded below as  $\min_{j \in S} \|\mathbf{u}_j^*\|_2 \geq g(\lambda_T)$ . Then  $\hat{\mathbf{u}}$  recovers the correct support, i.e.,  $S(\hat{\mathbf{u}}) = S(\mathbf{u}^*)$ .

**Theorem IV.3. ( $\ell_2$ -error bound)** Consider the linear model (4) satisfying assumptions (A1)–(A3). Let  $\lambda_T$  satisfies (10). Then with probability at least  $1 - \frac{2(m+1)}{(m_0+1)(m-m_0)}$ , the group lasso estimate satisfy

$$\left\| \begin{bmatrix} \mathbf{x}_0^* - \hat{\mathbf{x}}_0 \\ \mathbf{u}^* - \hat{\mathbf{u}} \end{bmatrix} \right\|_2 \leq \frac{2\sigma\sqrt{m_0+1}}{\sqrt{c_{\min}}} \left[ \sqrt{\frac{N}{T}} + \sqrt{\frac{\log(m_0+1)}{T}} \right] + \lambda_T \frac{\sqrt{m_0}}{c_{\min}}.$$

## V. SIMULATIONS

### VI. APPENDIX

**Additional Notation:** Let  $\mathbf{b} = [\mathbf{b}_1, \dots, \mathbf{b}_q] \in \mathbb{R}^{Mq}$ , where  $\mathbf{b}_j \in \mathbb{R}^M$ . For any set  $G \subseteq \{1, \dots, q\}$ , we let  $\mathbf{b}_G$  be the vector obtained by concatenating subvectors  $\mathbf{b}_j$  such that  $j \in G$ , and similarly, for  $\mathbf{b}_{G^c}$ . Let  $\mathcal{R}(\mathbf{b}_G) = \sum_{j \in G} \|\mathbf{b}_j\|_2$ . Note that with this notation, the penalty in (5) becomes  $\lambda_T \mathcal{R}(\mathbf{u})$ . It is a standard result in optimization theory that the dual norm of  $\mathcal{R}(\mathbf{b}_G)$ , denoted as  $\mathcal{R}^*(\mathbf{b}_G)$ , is

$$\mathcal{R}^*(\mathbf{b}_G) := \sup_{\mathcal{R}(\mathbf{c}) \leq 1} \mathbf{c}^\top \mathbf{b}_G = \max_{j \in G} \|\mathbf{b}_j\|_2. \quad (12)$$

*Proof.* The proof relies on the primal-dual witness (PDW) construction technique advocated by [?]. The PDW technique consists of following three steps<sup>2</sup>:

- (a) Set  $\hat{\mathbf{u}}_{S^c} = 0$ .
- (b) Determine  $(\hat{\mathbf{x}}_0, \hat{\mathbf{u}}_S)$  by solving the following oracle sub-problem:

$$\arg \min_{\substack{\mathbf{x} \in \mathbb{R}^n \\ \mathbf{u}_S \in \mathbb{R}^T}} \left\{ \frac{1}{2T} \|\mathbf{y} - \mathbf{O}\mathbf{x}_0 - \mathbf{J}_S \mathbf{u}_S\|_2^2 + \lambda_T \mathcal{R}(\mathbf{u}_S) \right\}. \quad (13)$$

Substituting  $\mathbf{x}_0 = \hat{\mathbf{x}}_0$  and  $\mathbf{u} = [\hat{\mathbf{u}}_S^\top \mathbf{0}^\top]^\top$  in (28), and solve for the vectors  $\hat{\mathbf{z}}_S \in \partial \mathcal{R}(\hat{\mathbf{u}}_S)$  and  $\hat{\mathbf{z}}_{S^c} \in \partial \mathcal{R}(\hat{\mathbf{u}}_{S^c})$ .

- (c) Check if the strict dual feasibility condition  $\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) < 1$  holds.

The PDW construction is said to be successful if the vector  $\hat{\mathbf{z}}_{S^c}$  satisfies the strict dual feasibility condition. Suppose that the PDW construction succeeds. In the view of Lemma ??,

<sup>2</sup>The PDW construction is not a numerical recipe for solving the group LASSO (5) problem, since solving the sub-problem in step (b) requires us to know about the *unknown* support set  $S$ . However, this proof technique helps us to solve the required result.

we observe that  $(\hat{\mathbf{x}}_0, \hat{\mathbf{u}}_S, \mathbf{0})$  obtained from steps (a) and (b) is the unique optimal solution of (5). Thus, all optimal input vectors are supported on the set  $S$ , i.e.,  $S(\hat{\mathbf{u}}) \subseteq S(\mathbf{u}^*)$ .

We now prove that  $\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) < 1$  with probability at least  $1 - 2/\bar{s}$ , where  $\bar{s} = |S^c| = m - m_0$ . We begin with fact that  $\hat{\mathbf{u}}_{S^c} = \hat{\mathbf{u}}_{S^c}^* = 0$ . Invoking this fact, along with the model (4), in (27a) and (28) results in the following block matrix form of KKT condition:

$$\begin{aligned} \frac{1}{T} \underbrace{\begin{bmatrix} \mathbf{O}^\top \mathbf{O} & \mathbf{O}^\top \mathbf{J}_S \\ \mathbf{J}_S^\top \mathbf{O} & \mathbf{J}_S^\top \mathbf{J}_S \end{bmatrix}}_{\mathbf{M}_S^\top \mathbf{M}_S} \underbrace{\begin{bmatrix} \mathbf{x}_0^* - \hat{\mathbf{x}}_0 \\ \mathbf{u}_S^* - \hat{\mathbf{u}}_S \end{bmatrix}}_{\Delta} + \frac{1}{T} \underbrace{\begin{bmatrix} \mathbf{O}^\top \\ \mathbf{J}_S^\top \end{bmatrix}}_{\mathbf{M}_S^\top} \mathbf{v} = \lambda_T \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{z}}_S \end{bmatrix}, \\ \frac{1}{T} \underbrace{\begin{bmatrix} \mathbf{J}_{S^c}^\top \mathbf{O} & \mathbf{J}_{S^c}^\top \mathbf{J} \end{bmatrix}}_{\mathbf{J}_{S^c}^\top \mathbf{M}_S} \Delta + \frac{1}{T} \mathbf{J}_{S^c}^\top \mathbf{v} = \lambda_T \hat{\mathbf{z}}_{S^c}. \end{aligned} \quad (14)$$

From (14) and the assumed invertibility condition (7),  $\Delta$  can be solved as

$$\Delta = -(\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \mathbf{M}_S^\top \mathbf{v} + \lambda_T T (\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{z}}_S \end{bmatrix}. \quad (15)$$

Substituting this equation into (14) and simplifying yields us

$$\begin{aligned} \hat{\mathbf{z}}_{S^c} = \underbrace{\mathbf{J}_{S^c}^\top \mathbf{M}_S (\mathbf{M}_S^\top \mathbf{M}_S)^{-1}}_{\mu} \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{z}}_S \end{bmatrix} \\ + \underbrace{\mathbf{J}_{S^c}^\top [\mathbf{I} - \mathbf{M}_S (\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \mathbf{M}_S^\top]}_{\eta_{S^c}} \left( \frac{\mathbf{v}}{\lambda_T T} \right). \end{aligned} \quad (16)$$

By the triangle inequality, it follows that  $\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) \leq \mathcal{R}^*(\mu) + \mathcal{R}^*(\eta_{S^c})$ . From the mutual incoherence condition (8), we deduce that  $\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) \leq \alpha + \mathcal{R}^*(\eta_{S^c})$ . Suppose that  $\mathcal{R}^*(\eta_{S^c}) \leq 0.5(1-\alpha)$ . Then, it is easy to see that  $\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) \leq 0.5(1+\alpha) < 1$ , which establishes the strict dual feasibility condition. Thus, it suffices to prove that  $\mathcal{R}^*(\eta_{S^c}) \leq 0.5(1-\alpha)$  holds with the claimed probability.

Define  $\Pi_{S^\perp} = \mathbf{I} - \mathbf{M}_S (\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \mathbf{M}_S^\top$ . Then,  $\eta_{S^c} = \mathbf{J}_{S^c}^\top \Pi_{S^\perp} (\mathbf{v}/\lambda_T T)$ , and from (12), we note that

$$\mathcal{R}^*(\eta_{S^c}) = \max_{j \in S^c} \|\mathbf{J}_j^\top \Pi_{S^\perp} (\mathbf{v}/\lambda_T T)\|_2. \quad (17)$$

Following similar steps in [], an elementary calculation shows that  $\|\mathbf{J}_j^\top \Pi_{S^\perp} (\mathbf{v}/\lambda_T T)\|_2$  is a Lipschitz with constant  $C/\lambda\sqrt{T}$ , where  $C$  is the group normalization bound defined in (A1). Hence, from Gaussian concentration of measure for Lipschitz functions [39], we have

$$\begin{aligned} \mathbb{P} \left( \left\| \mathbf{J}_j^\top \Pi_{S^\perp} \left( \frac{\mathbf{v}}{\lambda_T T} \right) \right\|_2 \geq \mathbb{E} \left[ \left\| \mathbf{J}_j^\top \Pi_{S^\perp} \left( \frac{\mathbf{v}}{\lambda_T T} \right) \right\|_2 \right] + \delta \right) \\ \leq 2 \exp \left( -\frac{T \lambda_T^2 \delta^2}{2C^2 \sigma^2} \right). \end{aligned}$$

Instead, using the Sudakov-Fernique comparison principle [], we can bound expectation term as

$$\mathbb{E} \left[ \left\| \mathbf{J}_j^\top \Pi_{S^\perp} \left( \frac{\mathbf{v}}{\lambda_T T} \right) \right\|_2 \right] \leq \frac{2C\sigma}{\lambda} \sqrt{\frac{N}{T}},$$

where  $N$  is the row dimension of  $\mathbf{J}_i^\top$ . From the above two inequalities, we may conclude that

$$\begin{aligned} \mathbb{P} \left( \left\| \mathbf{J}_j^\top \Pi_{S^\perp} \left( \frac{\mathbf{v}}{\lambda_T T} \right) \right\|_2 \geq 2 \frac{C\sigma}{\lambda} \sqrt{\frac{N}{T}} + \delta \right) \\ \leq 2 \exp \left( -\frac{T\lambda_T^2 \delta^2}{2C^2 \sigma^2} \right). \end{aligned} \quad (18)$$

Set  $\delta^2 = 4 \frac{\log \bar{s} C^2 \sigma^2}{\lambda_T^2 T}$ , where  $\bar{s} := m - m_0 = |S^c|$ . Then, from (10), we have  $(1 - \alpha) \geq \frac{2C\sigma}{\lambda_T} (\sqrt{N/T} + \sqrt{\log \bar{s}/T})$ . Thus,

$$\mathbb{P} \left( \left\| \mathbf{J}_j^\top \Pi_{S^\perp} \left( \frac{\mathbf{v}}{\lambda_T T} \right) \right\|_2 \geq 1 - \alpha \right) \leq 2 \exp(-2 \log \bar{s}).$$

Apply union bound over all indices in  $S^c$  to note that

$$\mathbb{P} \left( \max_{j \in S^c} \left\| \mathbf{J}_j^\top \Pi_{S^\perp} \left( \frac{\mathbf{v}}{\lambda_T T} \right) \right\|_2 \geq 1 - \alpha \right) \leq 2 \exp(-\log \bar{s}).$$

Finally, from (17), we have

$$\begin{aligned} \mathbb{P}[\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) \leq 0.5(1 - \alpha)] &\geq \mathbb{P}[\mathcal{R}^*(\hat{\mathbf{z}}_{S^c}) \leq (1 - \alpha)] \\ &\geq 1 - 2/\bar{s}, \end{aligned} \quad (19)$$

which establishes the required strict dual feasibility.

**Part 2:** We begin by bounding the dual norm of  $\mathbf{u}_S^* - \hat{\mathbf{u}}_S$ . Let  $\Gamma_j$  be as in the statement of Theorem IV.2, and  $\Delta$  be as in (14). Then, from (15), note the following:

$$\begin{aligned} \mathbf{u}_j^* - \hat{\mathbf{u}}_j &= \Gamma_j \Delta = \underbrace{-\Gamma_j (\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \mathbf{M}_S^\top \mathbf{v}}_{\mathbf{r}_{j,1}} \\ &\quad + \underbrace{\lambda_T T \Gamma_j (\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{z}}_S \end{bmatrix}}_{\mathbf{r}_{j,2}}. \end{aligned} \quad (20)$$

By the triangle inequality, we have  $\|\mathbf{u}_j^* - \hat{\mathbf{u}}_j\|_2 \leq \|\mathbf{r}_{j,1}\|_2 + \|\mathbf{r}_{j,2}\|_2$ . Taking maximum on both sides yields us

$$\max_{j \in S} \underbrace{\|\mathbf{u}_j^* - \hat{\mathbf{u}}_j\|_2}_{\mathcal{R}^*(\mathbf{u}_S^* - \hat{\mathbf{u}}_S)} \leq \underbrace{\max_{j \in S} \|\mathbf{r}_{j,1}\|_2}_{\mathcal{R}^*(\mathbf{r}_1)} + \lambda_T T \underbrace{\max_{j \in S} \|\mathbf{r}_{j,2}\|_2}_{\mathcal{R}^*(\mathbf{r}_2)}, \quad (21)$$

where  $\mathbf{r}_i, i \in \{1, 2\}$ , is formed by concatenating the vectors  $\mathbf{r}_{i,j}$  for  $j \in S$ . Let  $j^* = \arg \min_{j \in S} \|\mathbf{u}_j^*\|_2$ , and suppose that  $g(\lambda_T) \geq \mathcal{R}^*(\mathbf{r}_1) + \mathcal{R}^*(\mathbf{r}_2)$ . Consider the following:

$$\|\mathbf{u}_j^*\|_2 - \|\hat{\mathbf{u}}_j\|_2 \stackrel{(21)}{\leq} \mathcal{R}^*(\mathbf{r}_1) + \mathcal{R}^*(\mathbf{r}_2) \leq g(\lambda_T) < \|\mathbf{u}_{j^*}^*\|_2.$$

For the final inequality we invoked theorem hypothesis in part (2). Rearranging above inequality yields us  $\|\hat{\mathbf{u}}_j\|_2 \geq \|\mathbf{u}_j^*\|_2 - \|\mathbf{u}_{j^*}^*\|_2 > 0$ , for all  $j \in S$ . Thus, for the LASSO estimate (5),  $\hat{\mathbf{u}}_j \neq 0$  for all  $j \in S$ . Instead, in part (1) we showed that  $\hat{\mathbf{u}}_j = 0$  for all  $j \in S^c$ . Hence, it follows that  $S(\hat{\mathbf{u}}) = S(\mathbf{u}^*)$ .

We proceed to show that  $g(\lambda_T) \geq \mathcal{R}^*(\mathbf{r}_1) + \mathcal{R}^*(\mathbf{r}_2)$  with high probability. Notice that both  $\mathcal{R}^*(\mathbf{r}_1)$  and  $\mathcal{R}^*(\mathbf{r}_2)$  are random; however,  $\mathcal{R}^*(\mathbf{r}_2)$  is bounded above by a deterministic constant, almost surely. In fact,

$$\frac{\mathcal{R}^*(\mathbf{r}_2)}{\lambda_T T} = \max_{j \in S} \left\| \sum_{l \in S} \Gamma_j (\mathbf{M}^\top \mathbf{M})^{-1} \Gamma_l^\top \hat{\mathbf{z}}_l \right\|_2$$

$$\begin{aligned} &\leq \max_{j \in S} \sum_{l \in S} \left\| \Gamma_j (\mathbf{M}^\top \mathbf{M})^{-1} \Gamma_l^\top \right\|_2 \|\hat{\mathbf{z}}_l\|_2 \\ &\leq \underbrace{\max_{j \in S} \sum_{l \in S} \left\| \Gamma_j (\mathbf{M}^\top \mathbf{M})^{-1} \Gamma_l^\top \right\|_2}_{\xi_2}. \end{aligned} \quad (22)$$

The final inequality follows because  $\|\hat{\mathbf{z}}_l\|_2 \leq 1$ , for all  $l \in S$ ; follows from KKT conditions.

We now upper bound  $\mathcal{R}^*(\mathbf{r}_1) = \max_{j \in S} \|\mathbf{r}_{j,1}\|_2$ . From Lemma VI.1, for a  $j \in S$ , notice that  $\|\mathbf{r}_{j,1}\|_2$  is a Lipschitz with constant  $1/\sqrt{T c_{\min}}$ . Thus, once again, by Gaussian concentration of measure phenomenon [], we obtain

$$\mathbb{P}(\|\mathbf{r}_{j,1}\|_2 \geq \mathbb{E}[\|\mathbf{r}_{j,1}\|_2] + \gamma) \leq 2 \exp \left( -\frac{\gamma^2 T c_{\min}}{2\sigma^2} \right), \quad (23)$$

for any  $\gamma > 0$ . Instead, from the Sudakov-Fernique comparison principle [], it can be shown that  $\mathbb{E}[\|\mathbf{r}_{j,1}\|_2] \leq 2\sigma\sqrt{N/T c_{\min}}$ . Invoking this fact in (23) yields

$$\mathbb{P} \left( \|\mathbf{r}_{j,1}\|_2 \geq 2\sigma\sqrt{\frac{N}{T c_{\min}}} + \gamma \right) \leq 2 \exp \left( -\frac{\gamma^2 T c_{\min}}{2\sigma^2} \right).$$

Now, take union bound over all elements in  $S$ , and set  $\gamma^2 = \frac{4\sigma^2 \log m_0}{T c_{\min}}$ , where  $|S| = m_0$ , to conclude that

$$\mathbb{P} \left( \mathcal{R}^*(\mathbf{r}_2) \geq \frac{2\sigma}{\sqrt{c_{\min}}} \left( \sqrt{\frac{N}{T}} + \sqrt{\frac{\log m_0}{T}} \right) \right) \leq \frac{2}{m_0}. \quad (24)$$

Thus, from (11), (22), and (24), we can assert that

$$\mathbb{P}[\mathcal{R}^*(\mathbf{r}_1) + \mathcal{R}^*(\mathbf{r}_2) \leq g(\lambda_T)] \geq 1 - \frac{2}{m_0}$$

Finally, to obtain the probability claimed in Theorem IV.2, proceed as follows. Let  $G_1 = \{\mathcal{R}^*(\hat{\mathbf{z}}_S^c) \leq 0.5(1 - \alpha)\}$  and  $G_2 = \{\mathcal{R}^*(\mathbf{r}_1) + \mathcal{R}^*(\mathbf{r}_2) \leq g(\lambda_T)\}$  be two events. Then,

$$\mathbb{P}(G_1 \cap G_2) = 1 - \mathbb{P}(G_1^c \cup G_2^c) \geq 1 - [\mathbb{P}(G_1^c) + \mathbb{P}(G_2^c)].$$

Substituting (19) and (22) in the above inequality and simplifying yields the desired result.  $\square$

*Proof.* Suppose that the LASSO estimate (5) satisfies  $\hat{\mathbf{u}}_{S^c} = \mathbf{0}$ . Then, from the KKT conditions (27a) and (28) and the fact  $\mathbf{u}_{S^c}^* = \mathbf{0}$ , we note<sup>3</sup> that

$$\underbrace{\begin{bmatrix} \mathbf{x}_0^* - \hat{\mathbf{x}}_0 \\ \mathbf{u}_S^* - \hat{\mathbf{u}}_S \end{bmatrix}}_{\Delta} = - \underbrace{(\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \mathbf{M}_S^\top \mathbf{v}}_{\mathbf{r}_1} + \underbrace{\lambda_T T (\mathbf{M}_S^\top \mathbf{M}_S)^{-1} \begin{bmatrix} \mathbf{0} \\ \hat{\mathbf{z}}_S \end{bmatrix}}_{\mathbf{r}_2},$$

Furthermore,

$$\underbrace{\begin{bmatrix} \mathbf{x}_0^* - \hat{\mathbf{x}}_0 \\ \mathbf{u} - \hat{\mathbf{u}} \end{bmatrix}}_{\Delta_{\text{full}}} = \begin{bmatrix} \Delta \\ \mathbf{0} \end{bmatrix},$$

and, as a result,  $\|\Delta_{\text{full}}\|_2 = \|\Delta\|_2$ .

It now suffices to upper bound  $\|\Delta\|_2$ . First, observe that

$$\begin{aligned} \|\Delta\|_2 &\leq \|\mathbf{r}_1\|_2 + \|\mathbf{r}_2\|_2 \\ &\leq \|\mathbf{r}_1\|_2 + \lambda_T \left\| (\mathbf{M}_S^\top \mathbf{M}_S / T)^{-1} \right\|_2 \|\hat{\mathbf{z}}_S\|_2^2 \end{aligned}$$

<sup>3</sup>A similar computation is also used in prove part (1) of Theorem IV.2.

$$\leq \|\mathbf{r}_1\|_2 + \frac{\lambda_T}{c_{\min}} \sqrt{m_0} \quad (25)$$

For the final inequality, we invoked Assumption (A3) and the fact that  $\|\hat{\mathbf{z}}_S\|_2^2 = \sum_{j \in S} \|\hat{\mathbf{z}}_j\|_2^2 \leq |S| = m_0$ .

To bound  $\|\mathbf{r}_1\|_2$ , proceed as follows. Define  $\mathbf{u}_0 = \mathbf{x}_0$ . Using the arguments presented in the proof of part (2), Theorem IV.2, we can show that

$$\max_{j \in S \cup \{0\}} \|\mathbf{u}_j^* - \hat{\mathbf{u}}_j\|_2 \geq \underbrace{\frac{2\sigma}{\sqrt{c_{\min}}} \left( \sqrt{\frac{N}{T}} + \sqrt{\frac{\log(m_0 + 1)}{T}} \right)}_{\zeta}$$

holds with probability at most  $2/(m_0 + 1)$ . On the other hand, since  $\|\mathbf{r}_1\|_2 / \sqrt{m_0 + 1} \leq \max_{j \in S \cup \{0\}} \|\mathbf{u}_j^* - \hat{\mathbf{u}}_j\|_2$ , from the above inequality, we have  $\mathbb{P}(\|\mathbf{r}_1\|_2 \geq \zeta \sqrt{m_0 + 1}) \leq 2/(m_0 + 1)$ . Finally, from (25), we have

$$\mathbb{P}\left(\|\Delta\|_2 \leq \zeta \sqrt{m_0 + 1} + \frac{\lambda_T \sqrt{m_0}}{c_{\min}}\right) \geq 1 - \frac{2}{m_0 + 1}. \quad (26)$$

Note that the above bound is valid because we assumed that  $\hat{\mathbf{u}}_{S^c} = \mathbf{0}$ , i.e.,  $S(\hat{\mathbf{u}}) \subseteq S(\mathbf{u}^*)$ .

To obtain the probability claimed in the statement of Theorem IV.3, proceed as follows. Define the events  $G_1 = \{S(\hat{\mathbf{u}}) \subseteq S(\mathbf{u}^*)\}$  and  $G_2 = \{\|\Delta\|_2 \leq \zeta \sqrt{m_0 + 1} + \lambda_T \sqrt{m_0}/c_{\min}\}$ . Since  $\lambda_T$  satisfies (10), from (19), we have  $\mathbb{P}(G_1^c) \leq 2/(m - m_0)$ . Putting together the pieces, we can see that

$$\begin{aligned} \mathbb{P}(G_1 \cap G_2) &= 1 - \mathbb{P}(G_1^c \cup G_2^c) \geq 1 - [\mathbb{P}(G_1^c) + \mathbb{P}(G_2^c)] \\ &\geq 1 - \frac{2(m + 1)}{(m_0 + 1)(m - m_0)}. \end{aligned}$$

The proof is now complete.  $\square$

**Lemma VI.1.** The function  $\mathbf{v} \mapsto \|\Gamma_j(\mathbf{M}_S^T \mathbf{M}_S)^{-1} \mathbf{M}_S^T \mathbf{v}\|_2$  is a Lipschitz with constant  $1/\sqrt{T c_{\min}}$ .

*Proof.* Let  $\mathbf{X}_j = \Gamma_j(\mathbf{M}_S^T \mathbf{M}_S)^{-1} \mathbf{M}_S^T$ . For any pair of vectors  $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^T$ , we have

$$\begin{aligned} |\|\mathbf{X}_j \mathbf{v}_1\|_2 - \|\mathbf{X}_j \mathbf{v}_2\|_2| &\leq \|\mathbf{X}_j(\mathbf{v}_1 - \mathbf{v}_2)\|_2 \\ &\leq \|\mathbf{X}_j\|_2 \|\mathbf{v}_1 - \mathbf{v}_2\|_2. \end{aligned}$$

The above inequality establishes that  $\|\mathbf{X}_j \mathbf{v}\|_2$  is a Lipschitz with constant  $\|\mathbf{X}_j\|_2$ , which can be upper bounded as

$$\begin{aligned} \|\mathbf{X}_j\|_2^2 &= \lambda_{\max}(\mathbf{X}_j \mathbf{X}_j^T) \\ &= \lambda_{\max}(\Gamma_j(\mathbf{M}_S^T \mathbf{M}_S)^{-1} \mathbf{M}_S^T \mathbf{M}_S (\mathbf{M}_S^T \mathbf{M}_S)^{-1} \Gamma_j^T) \\ &= \lambda_{\max}(\Gamma_j(\mathbf{M}_S^T \mathbf{M}_S)^{-1} \Gamma_j^T) \\ &\leq \lambda_{\max}((\mathbf{M}_S^T \mathbf{M}_S)^{-1}) = 1/T c_{\min}. \end{aligned}$$

The inequality follows from Cauchy Interlace Theorem.  $\square$

**Proposition VI.2. (KKT Conditions)** A necessary and sufficient condition for  $[\mathbf{x}_0, \mathbf{u}_1, \dots, \mathbf{u}_m]$  to be solution of (5) is

$$-\frac{1}{T} \mathbf{O}^T [\mathbf{y} - \mathbf{O} \mathbf{x}_0 - \mathbf{J} \mathbf{u}] = \mathbf{0}, \quad (27a)$$

$$-\frac{1}{T} \mathbf{J}^T [\mathbf{y} - \mathbf{O} \mathbf{x}_0 - \mathbf{J} \mathbf{u}] + \lambda_T \mathbf{z} = \mathbf{0}, \quad (27b)$$

for all  $i \in \{1, \dots, m\}$ . Here,  $\mathbf{z} = [\mathbf{z}_1^T, \dots, \mathbf{z}_m^T]^T \in \partial \mathcal{R}(\mathbf{u})$ , where  $\partial \mathcal{R}(\mathbf{u})$  denotes the sub-gradient of  $\hat{\mathbf{u}}$ , and  $\mathbf{z}_i$  is the subgradient of  $\|\mathbf{u}_i\|_2$  evaluated at  $\mathbf{u}_i$ , which is given by

$$\mathbf{z}_i = \begin{cases} \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|_2}, & \text{if } \mathbf{u}_i = \mathbf{0} \\ \in \{\mathbf{q} : \|\mathbf{q}\|_2 \leq 1\}, & \text{if } \mathbf{u}_i \neq \mathbf{0} \end{cases}$$

*Proof.* See [].  $\square$

To prove Theorem IV.2, we need the following block matrix form of KKT condition (27b).

**Corollary VI.3. (Partitioned form of KKT conditions)** For any set  $G \subset \{1, \dots, m\}$ , let  $\mathbf{J} = [\mathbf{J}_G \quad \mathbf{J}_{G^c}]$  and  $\mathbf{u} = [\mathbf{u}_G \quad \mathbf{u}_{G^c}]$ . Then, the equation (27b) can be rewritten as

$$-\frac{1}{T} \begin{bmatrix} \mathbf{J}_G^T \\ \mathbf{J}_{G^c}^T \end{bmatrix} [\mathbf{y} - \mathbf{O} \mathbf{x}_0 - \mathbf{J}_G \mathbf{u}_G - \mathbf{J}_{G^c} \mathbf{u}_{G^c}] + \lambda_T \begin{bmatrix} \mathbf{z}_G \\ \mathbf{z}_{G^c} \end{bmatrix} = \mathbf{0}, \quad (28)$$

where  $\mathbf{z}_G$  is obtained by concatenating sub-vectors  $\mathbf{z}_j$ ,  $j \in G$ , and similarly, for  $\mathbf{z}_{G^c}$ .

## REFERENCES