

10-25-2014

Screening Rules for Overlapping Group Lasso

Seunghak Lee
Carnegie Mellon University

Eric P. Xing
Carnegie Mellon University, epxing@cs.cmu.edu

Follow this and additional works at: http://repository.cmu.edu/machine_learning

 Part of the [Theory and Algorithms Commons](#)

This Working Paper is brought to you for free and open access by the School of Computer Science at Research Showcase @ CMU. It has been accepted for inclusion in Machine Learning Department by an authorized administrator of Research Showcase @ CMU. For more information, please contact research-showcase@andrew.cmu.edu.

SCREENING RULES FOR OVERLAPPING GROUP LASSO

BY SEUNGHAK LEE AND ERIC P. XING

Carnegie Mellon University

Recently, to solve large-scale lasso and group lasso problems, screening rules have been developed, the goal of which is to reduce the problem size by efficiently discarding zero coefficients using simple rules independently of the others. However, screening for overlapping group lasso remains an open challenge because the overlaps between groups make it infeasible to test each group independently. In this paper, we develop screening rules for overlapping group lasso. To address the challenge arising from groups with overlaps, we take into account overlapping groups only if they are inclusive of the group being tested, and then we derive screening rules, adopting the dual polytope projection approach. This strategy allows us to screen each group independently of each other. In our experiments, we demonstrate the efficiency of our screening rules on various datasets.

1. Introduction. We propose efficient screening rules for regression with the overlapping group lasso penalty. Our goal is to develop simple rules to discard groups with zero coefficients in the optimization problem with the following form:

$$(1.1) \quad \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2,$$

where $\mathbf{X} \in \mathbb{R}^{N \times J}$ is the input data for J inputs and N samples, $\mathbf{y} \in \mathbb{R}^{N \times 1}$ is the output vector, $\boldsymbol{\beta} \in \mathbb{R}^{J \times 1}$ is the vector of regression coefficients, $n_{\mathbf{g}}$ is the size of group \mathbf{g} , and λ is a regularization parameter that determines the sparsity of $\boldsymbol{\beta}$. In this setting, \mathcal{G} represents a set of groups of coefficients, defined *a priori*, and we allow arbitrary overlap between different groups, hence “overlapping” group lasso. Overlapping group lasso is a general model that subsumes lasso (Tibshirani, 1996), group lasso (Yuan and Lin, 2006), sparse group lasso (Simon et al., 2013), composite absolute penalties (Zhao, Rocha and Yu, 2009), and tree lasso (Zhao, Rocha and Yu, 2009; Kim et al., 2012) with ℓ_1/ℓ_2 penalty because they are a specific form of overlapping group lasso.

In this paper, we do not consider the latent group lasso proposed by Jacob et al. (Jacob, Obozinski and Vert, 2009), where support is defined by the union of groups with nonzero coefficients. Instead, we consider the

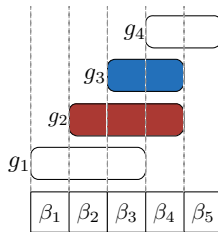


FIG 1. For screening test on group g_2 , we consider g_3 but disregard g_1 and g_4 to enable the independent screening test.

overlapping group lasso in the formulation of (1.1), where support is defined by the complement of the union of groups with zero coefficients (Jenatton, Audibert and Bach, 2011; Yuan, Liu and Ye, 2011). Therefore, unlike the latent group lasso, simple conversion from (1.1) to nonoverlapping group lasso problems by duplicating features overlapped between different groups is infeasible.

Recently, to solve (1.1) efficiently, fast algorithms have been developed (Yuan, Liu and Ye, 2011; Chen et al., 2012; Deng, Yin and Zhang, 2013) (we refer readers to (Bach et al., 2012) for review of optimization with sparsity-inducing penalties); however, in many applications such as genome-wide association studies (Lee and Xing, 2012; Yang et al., 2010), the number of features (or number of groups) can be very large. In such cases, fast optimization of (1.1) is challenging because it requires us to sweep over all coefficients/groups of coefficients many times until the objective converges. Furthermore, parallelization of existing sequential algorithms for speedup is nontrivial.

The past years have seen the emergence of screening techniques that can discard zero coefficients using simple rules in a single sweep over all coefficients. Examples include Sasvi rules (Liu et al., 2013), dual polytope projection (DPP) rules (Wang et al., 2013), dome tests (Xiang and Ramadge, 2012), sphere tests (Xiang, Xu and Ramadge, 2011), SAFE rules (Ghaoui, Viallon and Rabbani, 2012), and strong rules (Tibshirani et al., 2012). Among these, strong rules are not exact (true nonzero coefficients can be mistakenly discarded), whereas the other tests are exact. Furthermore, Bonnefoy et al. (Bonnefoy et al., 2014) recently developed an approach that merges screening approach with first-order optimization algorithms for lasso; however, to the best of our knowledge, none of the existing screening methods can be applied to (1.1) with overlapping groups. DPP and strong rules can be used for nonoverlapping group lasso, and the others are developed only for lasso.

In this paper, we develop exact screening rules for overlapping group lasso. The proposed screening rules can efficiently discard groups with zero coefficients by looking at each group independently. As a result, after screening, the number of groups that potentially include nonzero coefficients is small, and thus we can reduce the size of (1.1) by reformulating it using only the groups that survived. We then employ an optimization technique to solve the reduced problem. The resultant solution is optimal because the screening rules are exact in the sense that nonzero coefficients in a global optimal solution are never mistakenly discarded. The key idea behind our approach is to consider only groups that are inclusive of tested groups, while ignoring the other overlapping groups to perform independent screening tests. For example, in Figure 1, when performing a screening test on group g_2 , we take into account only g_3 because g_3 is a subset of g_2 , whereas g_1 and g_4 are non-inclusive of g_2 . The contributions of this paper are as follows:

1. We develop novel overlapping group lasso screening (OLS) and sparse overlapping group lasso screening (SOLS) rules. Sparse overlapping group lasso is a special case of overlapping group lasso, as formulated in (3.23).
2. We show that the screening rule for nonoverlapping group lasso via DPP (Wang et al., 2013) (group DPP or GDPP) is also exact when applied to overlapping group lasso.

In our experiments, we demonstrate that OLS, SOLS, and GDPP give us significant speed-up against a solver without screening. For example, OLS and SOLS achieved a $3.7\times$ and $3\times$ speed-up on PIE image dataset, compared to an overlapping group lasso solver without screening. Furthermore, OLS and SOLS are substantially more efficient to discard features with zero coefficients than GDPP under various experimental settings, confirming that the proposed algorithms are capable of using overlapping groups for screening.

Notation. We refer matrices to boldface and uppercase letters; vectors to boldface and lowercase letters; and scalars to regular lowercase letters. Columns are indexed by subscripts (e.g., \mathbf{x}_j is the j -th column vector of the matrix \mathbf{X}). We refer \mathbf{g} or \mathbf{h} to a group of coefficients, and $\mathbf{w}_{\mathbf{g}}$ represents a sub-vector of \mathbf{w} , indexed by \mathbf{g} .

2. Background: Screening Rules via DPP. Recently, Wang et al. proposed screening rules via DPP for nonoverlapping group lasso (Wang et al., 2013). The DPP screening rules are derived as follows: first, we find a dual form of group lasso and its Karush-Kuhn-Tucker (KKT) conditions. Then, using the KKT conditions and the relationship between primal and

dual solutions, we find screening rules to discard groups with zero coefficients; however, such screening rules involve a dual optimal solution, which is unknown. Thus, the DPP approach finds a range of vectors that includes a dual optimal solution, which is easy to obtain, and uses it instead of a dual optimal solution for screening rules. Here we review the DPP screening rule for nonoverlapping group lasso because it gives us with a vehicle to derive screening rules for overlapping group lasso.

For the nonoverlapping group lasso, screening rules via DPP can be obtained by the following procedure:

1. Find KKT conditions for a dual of the nonoverlapping group lasso.
2. Find the range of a dual optimal for the nonoverlapping group lasso.
3. Derive screening rules by using the range of a dual optimal and the KKT conditions.

The primal of the nonoverlapping group lasso is defined by

$$\min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}'} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2,$$

where \mathcal{G}' is a set of nonoverlapping groups, and $J = \sum_{\mathbf{g}} n_{\mathbf{g}}$. We first represent the nonoverlapping group lasso in a dual form:

$$(2.1) \quad \sup_{\boldsymbol{\theta}} \left\{ \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{\lambda^2}{2} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2^2 : \|\mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}\|_2 \leq \sqrt{n_{\mathbf{g}}}, \forall \mathbf{g} \in \mathcal{G}' \right\},$$

where $\boldsymbol{\theta} \in \mathbb{R}^{N \times 1}$ is a vector of dual variables. In (2.1), one can see that its dual optimal $\boldsymbol{\theta}^*$ is a vector ¹ that is closest to $\frac{\mathbf{y}}{\lambda}$ among the ones that satisfy all the constraints, and such a solution can be obtained by projecting $\frac{\mathbf{y}}{\lambda}$ onto the set of constraints $\mathbf{F}' \equiv \left\{ \|\mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}\|_2 \leq \sqrt{n_{\mathbf{g}}}, \forall \mathbf{g} \in \mathcal{G}' \right\}$. In other words, $\boldsymbol{\theta}^* = P_{\mathbf{F}'(\mathbf{y}/\lambda)}$, where $P_{\mathbf{F}'}$ denotes the projection operator onto \mathbf{F}' .

The KKT conditions of (2.1) (Wang et al., 2013) are given by

$$(2.2) \quad \begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta}^* + \lambda\boldsymbol{\theta}^*, \\ \mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}^* &= \begin{cases} \sqrt{n_{\mathbf{g}}} \frac{\boldsymbol{\beta}_{\mathbf{g}}^*}{\|\boldsymbol{\beta}_{\mathbf{g}}^*\|_2}, & \text{if } \boldsymbol{\beta}_{\mathbf{g}}^* \neq \mathbf{0}, \\ \sqrt{n_{\mathbf{g}}}\mathbf{u}, \|\mathbf{u}\|_2 \leq 1, & \text{if } \boldsymbol{\beta}_{\mathbf{g}}^* = \mathbf{0}. \end{cases} \end{aligned}$$

Based on (2.2), one can obtain a screening rule as follows: if $\|\mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}^*\|_2 < \sqrt{n_{\mathbf{g}}}$, then $\boldsymbol{\beta}_{\mathbf{g}}^* = \mathbf{0}$. However, it is still unusable because the dual optimal

¹For simple notation, we denote $\boldsymbol{\theta}^*$ by a dual optimal solution given λ . We use notation $\boldsymbol{\theta}^*(\lambda)$ when we refer to a specific λ .

$\boldsymbol{\theta}^*$ is unknown. To address this problem, we estimate a range of vectors, denoted by Θ , that contains $\boldsymbol{\theta}^*$ based on $\boldsymbol{\theta}^*(\lambda_0)$, where $\lambda_0 \neq \lambda$. Specifically, to estimate Θ , we use the fact that $P_{\mathbf{F}'}$ is continuous and nonexpansive. Finally, the DPP screening rule for nonoverlapping group lasso is formulated as follows: if $\sup_{\boldsymbol{\theta} \in \Theta} \|\mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}\|_2 < \sqrt{n_{\mathbf{g}}}$, then $\boldsymbol{\beta}_{\mathbf{g}}^* = \mathbf{0}$. By finding a closed-form solution for the left-hand side of the rule, i.e., $\sup_{\boldsymbol{\theta} \in \Theta} \|\mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}\|_2$, we can obtain screening rules for nonoverlapping group lasso.

3. Overlapping Group Lasso Screening. Now we develop overlapping group lasso screening rules. The first challenge is that the groups are not separable, making independent tests infeasible. Second, it is also unclear how to make use of overlapping groups in screening rules. Intuitively, coefficients are more likely to be zero if they are involved in more groups; thus incorporating overlapping groups into screening rules can help discard more features.

We address the first challenge by considering only a set of groups that is a subset of the tested group. Consider the example in Figure 1. Suppose we want to test if $\boldsymbol{\beta}_{\mathbf{g}_2} = \mathbf{0}$, $\mathbf{g}_2 = \{2, 3, 4\}$ given three other groups: $\mathbf{g}_1 = \{1, 2, 3\}$, $\mathbf{g}_3 = \{3, 4\}$, and $\mathbf{g}_4 = \{4, 5\}$. In such a case, we consider only \mathbf{g}_3 for the screening test on \mathbf{g}_2 because \mathbf{g}_3 is included in \mathbf{g}_2 , allowing us to test \mathbf{g}_2 independent of other groups; we ignore \mathbf{g}_1 and \mathbf{g}_4 in testing \mathbf{g}_2 because they involve coefficients not included in \mathbf{g}_2 , preventing us from testing the groups independently. For the second challenge, we derive new screening rules that can exploit the overlapping groups for minimizing the left-hand side of screening rules. In other words, they discard more features as the number of overlapping groups, inclusive of a tested group, increases.

In §3.1, we start with a condition for $\boldsymbol{\beta}_{\mathbf{g}}^* = \mathbf{0}$ for overlapping group lasso that contains a dual optimal. Based on this condition, we derive a screening rule by replacing a dual optimal with a range that includes the dual optimal in §3.2. Finally, in §3.3, we present efficient algorithms for overlapping group lasso screening based on the screening rules obtained in §3.2.

3.1. Screening Condition for Overlapping Group Lasso. Let us start with the dual form of overlapping group lasso (see appendix for derivation of the dual form):

$$(3.1) \quad \begin{aligned} & \sup_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{\lambda^2}{2} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2^2 \\ & \text{subject to } \mathbf{X}^T \boldsymbol{\theta} = \mathbf{v}, \end{aligned}$$

where $\boldsymbol{\theta}$ is a vector of dual variables, and \mathbf{v} is a subgradient of $\sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$ with respect to $\boldsymbol{\beta}$. A dual optimal $\boldsymbol{\theta}^*$ can be obtained by projecting $\frac{\mathbf{y}}{\lambda}$ onto

$\mathbf{F} \equiv (\mathbf{X}^T \boldsymbol{\theta} = \mathbf{v})$, denoted by $P_{\mathbf{F}}\left(\frac{\mathbf{y}}{\lambda}\right)$. We will use the “non-expansiveness” property of this projection operator to derive screening rules in §3.2.

Next we derive the KKT conditions for overlapping group lasso. Introducing $\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$, (1.1) can be written as

$$(3.2) \quad \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{z}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$$

subject to $\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$.

Then, a Lagrangian of (3.2) is

$$(3.3) \quad L(\boldsymbol{\beta}, \mathbf{z}, \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{z}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2 + \lambda \boldsymbol{\theta}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}),$$

and the KKT conditions of (3.3) are as follows:

$$(3.4) \quad 0 \in \frac{\partial L(\boldsymbol{\beta}^*, \mathbf{z}^*, \boldsymbol{\theta}^*)}{\partial \boldsymbol{\beta}_{\mathbf{g}}} = -\lambda \mathbf{X}_{\mathbf{g}}^T \boldsymbol{\theta}^* + \lambda \mathbf{v}_{\mathbf{g}},$$

$$(3.5) \quad 0 = \nabla_{\mathbf{z}} L(\boldsymbol{\beta}^*, \mathbf{z}^*, \boldsymbol{\theta}^*) = \mathbf{z}^* - \lambda \boldsymbol{\theta}^*,$$

$$(3.6) \quad 0 = \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\beta}^*, \mathbf{z}^*, \boldsymbol{\theta}^*) = \lambda (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^* - \mathbf{z}^*),$$

where $\mathbf{v}_{\mathbf{g}}$ is a subgradient of $\sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$ with respect to $\boldsymbol{\beta}_{\mathbf{g}}$. From (3.5) and (3.6), we obtain a bridge between the primal and dual solutions:

$$(3.7) \quad \lambda \boldsymbol{\theta}^* = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*.$$

Let us define two sets of groups that overlap with group \mathbf{g} as follows:

$$(3.8) \quad \bar{\mathcal{G}}_1 = \{\mathbf{h} : \mathbf{h} \in \mathcal{G} - \mathbf{g}, \mathbf{h} \subseteq \mathbf{g}, \mathbf{h} \cap \mathbf{g} \neq \emptyset\},$$

$$(3.9) \quad \bar{\mathcal{G}}_2 = \{\mathbf{h} : \mathbf{h} \in \mathcal{G} - \mathbf{g}, \mathbf{h} \not\subseteq \mathbf{g}, \mathbf{h} \cap \mathbf{g} \neq \emptyset\},$$

where $\bar{\mathcal{G}}_1$ and $\bar{\mathcal{G}}_2$ are sets of groups overlapping with \mathbf{g} , and $\bar{\mathcal{G}}_1$ includes the groups that are inclusive of \mathbf{g} . Then, we denote $\mathbf{v}_{\mathbf{g}} = \sqrt{n_{\mathbf{g}}} [\gamma_1, \dots, \gamma_{n_{\mathbf{g}}}]^T$, where γ_j is given by

$$(3.10) \quad \gamma_j = u_j + \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} w_j + \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_2} s_j,$$

where u_j is a subgradient of $\|\boldsymbol{\beta}_{\mathbf{g}}\|_2$ with respect to β_j ; w_j and s_j are subgradients of $\|\boldsymbol{\beta}_{\mathbf{h}}\|_2$ with respect to β_j , where \mathbf{h} belongs to $\bar{\mathcal{G}}_1$ and $\bar{\mathcal{G}}_2$, respectively. The definition of ℓ_2 norm subgradient, i.e., $\|\mathbf{u}_{\mathbf{g}}\|_2 \leq 1$, $\|\mathbf{w}_{\mathbf{h}}\|_2 \leq 1$,

and $\|\mathbf{s}_h\|_2 \leq 1$, gives us

$$(3.11) \quad \sqrt{\sum_{j \in \mathbf{g}} u_j^2} = \sqrt{\sum_{j \in \mathbf{g}} \left(\gamma_j - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} w_j - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_2} s_j \right)^2} \leq 1,$$

where the equality holds when $\beta_{\mathbf{g}}^* \neq \mathbf{0}$. Plugging (3.4) into (3.11), a sufficient condition for $\beta_{\mathbf{g}}^* = \mathbf{0}$ is given by

$$(3.12) \quad \min_{\mathbf{w}, \mathbf{s}: \|\mathbf{w}_h\|_2 \leq 1, \|\mathbf{s}_h\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \boldsymbol{\theta}^* - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_h} w_j - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_2} \sqrt{n_h} s_j \right)^2} < \sqrt{n_{\mathbf{g}}}.$$

To screen each group independently (i.e., test using only the coefficients in group \mathbf{g}), we set $\mathbf{s}_h = \mathbf{0}$, for all $\mathbf{h} \in \bar{\mathcal{G}}_2$. This is a valid subgradient because it always satisfies $\|\mathbf{s}_h\|_2 \leq 1$. Given subgradients of the groups inclusive of \mathbf{g} , we have the following screening condition for $\beta_{\mathbf{g}}^* = \mathbf{0}$: if $b_{\mathbf{g}} < \sqrt{n_{\mathbf{g}}}$, then $\beta_{\mathbf{g}}^* = \mathbf{0}$, where $b_{\mathbf{g}}$ is defined by

$$(3.13) \quad b_{\mathbf{g}} \equiv \min_{\mathbf{w}: \|\mathbf{w}_h\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \boldsymbol{\theta}^* - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_h} w_j \right)^2}.$$

Note that $b_{\mathbf{g}}$ is an upper bound on the left-hand side of (3.12) due to the fixed $\mathbf{s}_h = \mathbf{0}$, for all $\mathbf{h} \in \bar{\mathcal{G}}_2$. If $b_{\mathbf{g}} < \sqrt{n_{\mathbf{g}}}$, then (3.12) holds, and thus $\beta_{\mathbf{g}}^* = \mathbf{0}$.

3.2. Screening Rules for Overlapping Group Lasso. So far, we have derived a condition for $\beta_{\mathbf{g}}^* = \mathbf{0}$; however, it is not yet usable for screening because $\boldsymbol{\theta}^*$ is unknown. Thus, by following the DPP approach by Wang et al. (Wang et al., 2013), we first find a region Θ that contains $\boldsymbol{\theta}^*$.

In (3.1), an optimal $\boldsymbol{\theta}^*$ is the projection of $\frac{\mathbf{y}}{\lambda}$ onto the constraint \mathbf{F} :

$$\boldsymbol{\theta}^* = P_{\mathbf{F}} \left(\frac{\mathbf{y}}{\lambda} \right) = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbf{F}} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2,$$

where $P_{\mathbf{F}}$ is the projection operator with \mathbf{F} which is a nonempty, closed convex subset of a Hilbert space (\mathbf{F} is nonempty because $\mathbf{0} \in \mathbf{F}$, and closed convex because it is an intersection of closed half-spaces). Thus, we can use the “non-expansiveness” property of $P_{\mathbf{F}}$ (Bertsekas et al., 2003), given by

$$(3.14) \quad \left\| P_{\mathbf{F}} \left(\frac{\mathbf{y}}{\lambda} \right) - P_{\mathbf{F}} \left(\frac{\mathbf{y}}{\lambda_0} \right) \right\|_2 = \|\boldsymbol{\theta}^*(\lambda) - \boldsymbol{\theta}^*(\lambda_0)\|_2 \leq \left\| \frac{\mathbf{y}}{\lambda} - \frac{\mathbf{y}}{\lambda_0} \right\|_2,$$

where λ_0 is a tuning parameter ($\lambda_0 > \lambda$), and $\boldsymbol{\theta}^*(\lambda) \equiv \boldsymbol{\theta}^*$, and $\boldsymbol{\theta}^*(\lambda_0)$ is a dual optimal solution given λ_0 . Here (3.14) shows that $\boldsymbol{\theta}^*(\lambda)$ lies within a sphere Θ centered at $\boldsymbol{\theta}^*(\lambda_0)$ with a radius of $\rho = \left\| \frac{\mathbf{y}}{\lambda} - \frac{\mathbf{y}}{\lambda_0} \right\|_2$. Based on this, we can represent $\boldsymbol{\theta}^*(\lambda) = \boldsymbol{\theta}^*(\lambda_0) + \mathbf{r}$, where $\|\mathbf{r}\|_2 \leq \rho$. By plugging it into (3.13) and maximizing the objective over \mathbf{r} , we have the following screening rule: if $b'_g < \sqrt{n_g}$, then $\boldsymbol{\beta}^* = \mathbf{0}$, where b'_g is defined by

$$(3.15) \quad b'_g \equiv \sup_{\mathbf{r}: \|\mathbf{r}\|_2 \leq \rho} \min_{\mathbf{w}: \|\mathbf{w}_h\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \{ \boldsymbol{\theta}^*(\lambda_0) + \mathbf{r} \} - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_h} w_j \right)^2}.$$

Notice that b'_g is an upper bound on b_g , and thus $b'_g < \sqrt{n_g} \Rightarrow b_g < \sqrt{n_g} \Rightarrow \boldsymbol{\beta}^* = \mathbf{0}$. With a little bit of algebra, we get our screening rule for overlapping group lasso.

THEOREM 1. *For the overlapping lasso problem, suppose that we are given an optimal dual solution $\boldsymbol{\theta}^*(\lambda_0)$. Then for $\lambda < \lambda_0$, $\boldsymbol{\beta}_g^*(\lambda) = \mathbf{0}$ if*

$$(3.16) \quad \min_{\mathbf{w}_h: \|\mathbf{w}_h\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_h} w_j \right)^2} < \sqrt{n_g} - \|\mathbf{X}_g\|_F \|\mathbf{y}\|_2 \left| \frac{1}{\lambda} - \frac{1}{\lambda_0} \right|.$$

PROOF. See Appendix B. □

Using the bridge between the primal and dual in (3.7), we can also obtain a screening rule in a primal form.

THEOREM 2. *For the overlapping lasso problem, suppose that we are given an optimal solution $\boldsymbol{\beta}^*(\lambda_0)$. Then for $\lambda < \lambda_0$, $\boldsymbol{\beta}_g^*(\lambda) = \mathbf{0}$ if*

$$(3.17) \quad \min_{\mathbf{w}_h: \|\mathbf{w}_h\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \frac{\mathbf{y} - \mathbf{X} \boldsymbol{\beta}^*(\lambda_0)}{\lambda_0} - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_h} w_j \right)^2} < \sqrt{n_g} - \|\mathbf{X}_g\|_F \|\mathbf{y}\|_2 \left| \frac{1}{\lambda} - \frac{1}{\lambda_0} \right|.$$

We also note that Theorem 2 can be employed to solve lasso problems following a λ path $\{\lambda_1, \lambda_2, \dots, \lambda_T\}$ in a descending order. The λ path is determined *a priori*, and one may choose linearly, geometrically, or logarithmically spaced λ values. In the sequential version of screening, we first perform screening with λ_1 using λ' with $\boldsymbol{\beta}^*(\lambda') = \mathbf{0}$. We then run a solver using the remaining coefficients with their corresponding groups; its results become $\boldsymbol{\beta}^*(\lambda_1)$. Now, using λ_1 with $\boldsymbol{\beta}^*(\lambda_1)$, we perform screening for λ_2 . We repeat the above procedure for all remaining λ parameters. The following theorem shows sequential screening rule, where a screening rule for λ_t is constructed based on $\boldsymbol{\beta}^*(\lambda_{t-1})$, $t \geq 2$.

THEOREM 3. *For the overlapping lasso problem with a λ path $\{\lambda_1, \dots, \lambda_T\}$, $\lambda_{t-1} > \lambda_t$, $t = 2, \dots, T$, suppose that we are given an optimal solution $\beta^*(\lambda_{t-1})$. Then, $\beta_{\mathbf{g}}^*(\lambda_t) = \mathbf{0}$ if*

$$(3.18) \quad \min_{\mathbf{w}_{\mathbf{h}}: \|\mathbf{w}_{\mathbf{h}}\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \frac{\mathbf{y} - \mathbf{X}\beta^*(\lambda_{t-1})}{\lambda_{t-1}} - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_{\mathbf{h}}} w_j \right)^2} < \sqrt{n_{\mathbf{g}}} - \|\mathbf{X}_{\mathbf{g}}\|_F \|\mathbf{y}\|_2 \left| \frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right|.$$

We omit the proofs for Theorem 2 and Theorem 3 because it is straightforward to derive them from Theorem 1.

3.3. Screening Algorithms for Overlapping Group Lasso. To use Theorems 1, 2, 3 we need an efficient way to obtain the left-hand side. Instead of solving the left-hand side directly, we minimize an upper bound on the left-hand side because it can be quickly solved. Any upper bounds give us a valid screening rule, but the tighter the bound, the better the screening efficiency (it discards more features). Note that the goal of screening is to speed up optimization, and thus we intend to present a simple yet efficient algorithm.

We first present our algorithm and then verify that it minimizes an upper bound on the left-hand side. We adopt a simple coordinate descent-type approach, where each group is used for minimization one at a time. Suppose that we perform screening on group \mathbf{g} . We start with making two variables: $l \leftarrow 0$, $\mathbf{a} \leftarrow \mathbf{g}$, and a set of overlapping groups $\{\mathbf{h}_1, \dots, \mathbf{h}_K : \mathbf{h}_k \in \bar{\mathcal{G}}_1, k = 1, \dots, K\}$ (any ordering works for our purpose). For each group \mathbf{h}_k , we take the intersection between \mathbf{h}_k and \mathbf{a} , i.e., $\mathbf{h}'_k \leftarrow \mathbf{h}_k \cap \mathbf{a}$. If $\mathbf{h}'_k = \emptyset$, we skip \mathbf{h}_k and proceed to the next \mathbf{h}_{k+1} . If $\mathbf{h}'_k \neq \emptyset$, we take the following procedure. If $\|\mathbf{X}_{\mathbf{h}'_k} \boldsymbol{\theta}^*(\lambda_0)\|_2 \leq \sqrt{n_{\mathbf{h}'_k}}$, we set $l \leftarrow l$; otherwise $l \leftarrow l + z$, where $z = \sum_{j \in \mathbf{h}'_k} \left\{ \mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sqrt{n_{\mathbf{h}'_k}} w_j \right\}^2$, and $\{w_j : j \in \mathbf{h}'_k\}$ is determined by the following algorithm.

1. Set $d = 1$.
2. For each $j \in \mathbf{h}'_k$, we compute

$$(3.19) \quad w_j = \begin{cases} \frac{\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0)}{\sqrt{n_{\mathbf{h}'_k}}} & , \text{ if } \left| \frac{\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0)}{\sqrt{n_{\mathbf{h}'_k}}} \right| \leq \sqrt{d} \\ \text{sign} \left\{ \mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) \right\} \sqrt{d} & , \text{ otherwise,} \end{cases}$$

and update $d \leftarrow d - w_j^2$.

Algorithm 1: Screening for overlapping group Lasso

Input: $\mathbf{X}, \mathbf{y}, \lambda_{t-1}, \lambda_t$ ($\lambda_{t-1} > \lambda_t$), $\mathcal{G}, \boldsymbol{\theta}^*(\lambda_{t-1})$ ($= \frac{\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*(\lambda_{t-1})}{\lambda_{t-1}}$)

Output: \mathcal{T} (a set of groups with potential nonzero coefficients)

```

1  $\mathcal{T} \leftarrow \emptyset;$ 
2 for  $\mathbf{g} \in \mathcal{G}$  do
3    $\bar{\mathcal{G}}_1 = \{\mathbf{h} : \mathbf{h} \in \mathcal{G} - \mathbf{g}, \mathbf{h} \subseteq \mathbf{g}, \mathbf{h} \cap \mathbf{g} \neq \emptyset\};$ 
4    $\mathbf{a} \leftarrow \mathbf{g};$ 
5    $l \leftarrow 0;$ 
6   for  $\mathbf{h} \in \bar{\mathcal{G}}_1$  do
7      $\mathbf{h}' \leftarrow \mathbf{h} \cap \mathbf{a};$ 
8     if  $\|\mathbf{X}_{\mathbf{h}}^T \boldsymbol{\theta}^*(\lambda_{t-1})\|_2 > \sqrt{n_{\mathbf{h}}}$  then
9        $d \leftarrow 1;$ 
10      for  $j \in \mathbf{h}'$  do
11        if  $|\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1})| \leq \sqrt{dn_{\mathbf{h}}}$  then
12           $d \leftarrow d - \left(\frac{\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1})}{\sqrt{n_{\mathbf{h}}}}\right)^2;$ 
13        else
14           $l \leftarrow l + \left[\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1}) - \sqrt{dn_{\mathbf{h}}} \text{sign}\{\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1})\}\right]^2;$ 
15          break;
16       $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{h}';$ 
17   if  $\sqrt{l + \|\mathbf{X}_{\mathbf{a}}^T \boldsymbol{\theta}^*(\lambda_{t-1})\|_2^2} \geq \sqrt{n_{\mathbf{g}}} - \|\mathbf{X}_{\mathbf{g}}\|_F \|\mathbf{y}\|_2 \left|\frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}}\right|$  then
18      $\mathcal{T} \leftarrow \mathcal{T} \cup \mathbf{g};$ 

```

We then set $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{h}'_k$, and iterate this procedure over all groups in $\bar{\mathcal{G}}_1$.

Finally, a minimized left-hand side for the screening rules is $\sqrt{l + \|\mathbf{x}_{\mathbf{a}}^T \boldsymbol{\theta}^*(\lambda_0)\|_2^2}$. This algorithm in a sequential setting for overlapping group lasso (screening for λ_t given λ_{t-1} , and $\lambda_t \in \{\lambda_1, \dots, \lambda_T\}$) is summarized in Algorithm 1.

Now, we show that this algorithm minimizes an upper bound on the left-hand side of Theorem 1. The key idea is that, at the $(k+1)$ -th iteration, we minimize an upper bound on the bound obtained at the k -th iteration. We denote \mathbf{a} by the set of coefficients to be processed, and k by the iteration counter, initialized by $\mathbf{a} \leftarrow \mathbf{g}$ and $k = 1$. At the k -th iteration, the squared left-hand side is bounded as follows:

$$(3.20) \quad \min_{\mathbf{w}_{\mathbf{h}}, \forall \mathbf{h} \in \bar{\mathcal{G}}_1} \left[\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_{\mathbf{h}}} w_j \right)^2 \right]$$

$$(3.21) \quad \leq \min_{\mathbf{w}_{\mathbf{h}}, \forall \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k} \left[\sum_{j \in \mathbf{g} - \mathbf{h}_k} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k} \sqrt{n_{\mathbf{h}}} w_j \right)^2 \right] \\ + \min_{\mathbf{w}_{\mathbf{h}_k}} \left\| \mathbf{X}_{\mathbf{h}_k}^T \boldsymbol{\theta}^*(\lambda_0) - \sqrt{n_{\mathbf{h}_k}} \mathbf{w}_{\mathbf{h}_k} \right\|_2^2.$$

To obtain the upper bound in (3.21), we set $w_j = 0$ for all $\{j : j \in \mathbf{h} \cap \mathbf{h}_k, \mathbf{h} \in$

$\bar{\mathcal{G}}_1 - \mathbf{h}_k$. Let us fix $\{\mathbf{w}_{\mathbf{h}} : \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k\}$ and then find the bound in (3.21). Since the first term is a constant due to fixed $\{\mathbf{w}_{\mathbf{h}}\}$, we find z that bounds the second term:

$$\min_{\mathbf{w}_{\mathbf{h}_k}} \left\| \mathbf{X}_{\mathbf{h}_k}^T \boldsymbol{\theta}^*(\lambda_0) - \sqrt{n_{\mathbf{h}_k}} \mathbf{w}_{\mathbf{h}_k} \right\|_2^2 \leq z.$$

Based on the subgradient condition $\|\mathbf{w}_{\mathbf{h}_k}\|_2 \leq 1$, if $\left\| \mathbf{X}_{\mathbf{h}_k}^T \boldsymbol{\theta}^*(\lambda_0) \right\|_2 \leq \sqrt{n_{\mathbf{h}_k}}$, we set $\left\| \mathbf{X}_{\mathbf{h}_k}^T \boldsymbol{\theta}^*(\lambda_0) - \sqrt{n_{\mathbf{h}_k}} \mathbf{w}_{\mathbf{h}_k} \right\|_2^2 = 0 = z$; otherwise, we find an upper bound z using the simple coordinate descent-type procedure with (3.19). It is easy to see that the procedure with (3.19) satisfies the subgradient condition $\|\mathbf{w}_{\mathbf{h}_k}\|_2 \leq 1$, and thus z is a valid upper bound. Then, we set $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{h}_k$ because \mathbf{h}_k is processed, and set $l \leftarrow z$.

Subsequently, we get an upper bound on (3.21):

$$\begin{aligned} & \min_{\mathbf{w}^{\mathbf{h}}, \forall \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k} \left[\sum_{j \in \mathbf{g} - \mathbf{h}_k} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k} \sqrt{n_{\mathbf{h}}} w_j \right)^2 \right] + l \\ (3.22) \quad & \leq \min_{\mathbf{w}^{\mathbf{h}}, \forall \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k - \mathbf{h}_{k+1}} \left[\sum_{j \in \mathbf{g} - \mathbf{h}_k - \mathbf{h}_{k+1}} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k - \mathbf{h}_{k+1}} \sqrt{n_{\mathbf{h}}} w_j \right)^2 \right] \\ & \quad + \min_{\mathbf{w}_{\mathbf{h}'_{k+1}}} \left\| \mathbf{X}_{\mathbf{h}'_{k+1}}^T \boldsymbol{\theta}^*(\lambda_0) - \sqrt{n_{\mathbf{h}'_{k+1}}} \mathbf{w}_{\mathbf{h}'_{k+1}} \right\|_2^2 + l, \end{aligned}$$

where $\mathbf{h}'_{k+1} = \mathbf{h}_{k+1} \cap \mathbf{a}$. Fixing $\{\mathbf{w}_{\mathbf{h}} : \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k - \mathbf{h}_{k+1}\}$, and setting $w_j = 0, \forall j \in \mathbf{h} \cap \mathbf{h}_k, \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_k - \mathbf{h}_{k+1}$, we minimize an upper bound z on $\left\| \mathbf{X}_{\mathbf{h}'_{k+1}}^T \boldsymbol{\theta}^*(\lambda_0) - \sqrt{n_{\mathbf{h}'_{k+1}}} \mathbf{w}_{\mathbf{h}'_{k+1}} \right\|_2^2$ using the procedure with (3.19), and $l \leftarrow l + z$. We iterate this procedure over all groups in $\bar{\mathcal{G}}_1$, i.e., $\{\mathbf{h} \in \bar{\mathcal{G}}_1\}$, resulting in an upper bound on the left-hand side as follows:

$$\min_{\mathbf{w}^{\mathbf{h}}, \forall \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_1 - \dots - \mathbf{h}_K} \left[\sum_{j \in \mathbf{g} - \mathbf{h}_1 - \dots - \mathbf{h}_K} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_1 - \dots - \mathbf{h}_K} \sqrt{n_{\mathbf{h}}} w_j \right)^2 \right] + l.$$

By setting $w_j = 0$ for all $j \in \mathbf{h} \cap \mathbf{h}_k, \mathbf{h} \in \bar{\mathcal{G}}_1 - \mathbf{h}_1 - \dots - \mathbf{h}_K$, we get an upper bound on the squared left-hand side, i.e., $l + \left\| \mathbf{x}_{\mathbf{a}}^T \boldsymbol{\theta}^*(\lambda_0) \right\|_2^2$; taking the square root on it, we obtain the left-hand side of Theorem 1.

We note that the DPP screening rule for nonoverlapping group lasso (GDPP) (Wang et al., 2013) is a special case of the proposed screening rules for overlapping group lasso. In Theorem 1, by setting $w_j = 0$ for all $j \in \mathbf{h}$, we obtain GDPP, where its left-hand side is an upper bound on that of our screening rules. This implies that GDPP is also an exact screening

Algorithm 2: Screening for sparse overlapping group lasso

Input: $\mathbf{X}, \mathbf{y}, \lambda_{t-1}, \lambda_t$ ($\lambda_{t-1} > \lambda_t$), $\mathcal{G}, \boldsymbol{\theta}^*(\lambda_{t-1})$ ($= \frac{\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*(\lambda_{t-1})}{\lambda_{t-1}}$)

Output: \mathcal{T} (a set of groups with potential nonzero coefficients)

```

1  $\mathcal{T} \leftarrow \emptyset;$ 
2 for  $\mathbf{g} \in \mathcal{G}$  do
3   if  $|\mathcal{G}| \geq 2$  then
4      $l \leftarrow 0;$ 
5     for  $j \in \mathcal{G}$  do
6       if  $|\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1})| > 1$  then
7          $l \leftarrow l + \left\{ \mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1}) - \text{sign} \left( \mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1}) \right) \right\}^2;$ 
8       else
9          $l \leftarrow \left\{ \mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_{t-1}) \right\}^2;$ 
10    if  $\sqrt{l} \geq \sqrt{n_{\mathbf{g}}} - \|\mathbf{X}_{\mathbf{g}}\|_F \|\mathbf{y}\|_2 \left| \frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right|$  then
11       $\mathcal{T} \leftarrow \mathcal{T} \cup \mathbf{g};$ 

```

rule for overlapping group lasso; however, GDPP would not be as efficient as Theorem 1 due to the lack of degrees of freedom to decrease its left-hand side. It is surprising that GDPP is applicable to overlapping group lasso because GDPP is derived under the assumption that groups do not overlap.

In practice, finding $\bar{\mathcal{G}}_1$ can be an algorithmic bottleneck (line 3 in Algorithm 1). To search for $\bar{\mathcal{G}}_1$ efficiently, we used a simple algorithm. We first sort the groups based on the smallest index of each group, resulting in $\mathcal{G} = \{\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(M)}\}$. To perform a screening test on $\mathbf{g}^{(m)}$, we find its $\bar{\mathcal{G}}_1$ by searching for the groups between $\mathbf{g}^{(m+1)}$ and $\mathbf{g}^{(m+W)}$, where W is the user-defined window size. With larger W , we may discard more features, but the computational complexity for screening increases linearly in W .

3.3.1. *Screening Algorithm for Sparse Overlapping Group Lasso.* Sparse overlapping group lasso is a special case of overlapping group lasso that includes ℓ_1 penalty, defined by

$$(3.23) \quad \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2.$$

For this model, we provide a simple and fast algorithm by considering only the individual coefficients in ℓ_1 penalty for $\bar{\mathcal{G}}_1$. Note that individual features can be considered as groups of size one, inclusive of other groups. Substituting $\bar{\mathcal{G}}_1$ in line 3 in Algorithm 1 by $\bar{\mathcal{G}}_1 = \{j : j \in \mathbf{g}\}$, we obtain a screening algorithm for sparse overlapping group lasso, summarized in Algorithm 2.

Algorithm 3: Finding a small λ that sets all coefficients to zero

Input: $\mathbf{X}, \mathbf{y}, \mathcal{G}, r$ ($0 < r < 1$ is a common ratio for a geometric series)
Output: λ' that discards all features

```

1  $\lambda_1 = \max_{\mathbf{g} \in \mathcal{G}} \frac{1}{\sqrt{n_{\mathbf{g}}}} \|\mathbf{X}_{\mathbf{g}}^T \mathbf{y}\|_2;$ 
2 for  $t = 2$  to  $T$  do
3    $\lambda_t = \lambda_{t-1} r;$ 
4    $\boldsymbol{\theta}^*(\lambda_{t-1}) = \frac{\mathbf{y}}{\lambda_{t-1}};$ 
5    $\mathcal{T} \leftarrow \text{Algorithm 1}(\mathbf{X}, \mathbf{y}, \lambda_t, \lambda_{t-1}, \mathcal{G}, \boldsymbol{\theta}^*(\lambda_{t-1}));$ 
6   if  $\mathcal{T} \neq \emptyset$  then
7      $\lambda' = \lambda_{t-1};$ 
8     break;

```

It is worthwhile to mention that Algorithm 2 is significantly faster than Algorithm 1 due to the lack of set operations in Algorithm 1. Furthermore, in our experiments, we observed that Algorithm 2 discards similar numbers of features to Algorithm 1 on various datasets. Thus, if the model in (3.23) is utilized in an application, Algorithm 2 would be appealing in terms of both its screening rejection power and speed.

Remark. For nonoverlapping group lasso, we find λ_{max} (the smallest λ that sets $\boldsymbol{\beta}^* = \mathbf{0}$) as follows: $\lambda_{max} = \max_{\mathbf{g} \in \mathcal{G}} \frac{1}{\sqrt{n_{\mathbf{g}}}} \|\mathbf{X}_{\mathbf{g}}^T \mathbf{y}\|_2$. However, this is not necessarily the smallest λ for zero solutions for overlapping group lasso. In fact, it is nontrivial to compute λ_{max} for overlapping group lasso because different groups are coupled through overlaps, preventing us from using the simple equation above. Instead, using a screening algorithm, we can find a small λ that sets all coefficients to zero, denoted by λ' . The key idea is that we decrease λ following a sequence until all coefficients are discarded by a screening algorithm. We denote λ' by the smallest λ that sets $\boldsymbol{\beta}^* = \mathbf{0}$ in our regularization path. This technique is summarized in Algorithm 3 with a geometric sequence of λ parameters.

4. Experiments. We demonstrate the efficiency of the proposed screening algorithms in terms of the screening rejection ratio and the speed. The rejection ratio is defined by the ratio of discarded coefficients to true zero coefficients, obtained by an overlapping group lasso solver without screening. Here we refer Algorithm 1 to OLS and Algorithm 2 to SOLS.

To the best of our knowledge, there are no existing screening rules for overlapping group lasso; however, we showed that GDPP can also be used for overlapping group lasso as an exact screening rule. Therefore, as a baseline screening algorithm, we used GDPP. Comparing OLS² and SOLS against

²We set the window size to 50 to search for the overlapping groups in $\bar{\mathcal{G}}_1$.

GDPP, we investigate the benefits of using overlapping groups for screening because GDPP makes no use of overlapping groups.

We used the following three image datasets ³ and one genome dataset for our experiments, chosen to cover the cases where $N \gg J$, $N \ll J$, and $N \approx J$, and a real-world application in genetics. **a)** The PIE image database (Sim, Baker and Bsat, 2002), which contains 11,554 face images of 68 people under different poses, illumination conditions, and expressions. The images are represented by $\mathbf{D} \in \mathbb{R}^{11554 \times 1024}$. We generated the response vector $\mathbf{y} \in \mathbb{R}^{11554 \times 1}$ by randomly choosing a feature in \mathbf{D} , and the rest of the features are concatenated to be the design matrix $\mathbf{X} \in \mathbb{R}^{11554 \times 1023}$. **b)** The Alzheimer’s disease (AD) dataset (Zhang et al., 2013), which contains 541 AD individuals, represented by 511,997 single nucleotide polymorphisms (the most common genetic variants). For the same individuals, the AD dataset contains 40,638 expression levels for known and predicted genes, splice variants, miRNAs, and non-coding RNA sequences in the brain region of the cerebellum. We used the genetic information for $\mathbf{X} \in \mathbb{R}^{541 \times 511997}$, and randomly choose a gene expression for $\mathbf{y} \in \mathbb{R}^{541 \times 1}$. This is a typical experimental setting for expression quantitative mapping (Lee and Xing, 2012), where the goal is to identify genetic variants that affect gene expression levels. **c)** The COIL database (Nayar, Nene and Murase, 1996) contains color images of 100 objects. For each object, 72 images are taken with different angles. We selected object 10, whose image data are represented by $\mathbf{D} \in \mathbb{R}^{72 \times 49152}$. \mathbf{X} and \mathbf{y} are generated in the same way as the PIE dataset. **d)** The digit recognition data (GISETTE) (Guyon et al., 2004), which contains 6,000 digit images of four and nine. Each sample contains 5,000 features, where 70% of features are constructed from MNIST dataset (LeCun et al., 1998), and 30% of them are artificially generated following a distribution of true features, resulting in $\mathbf{X} \in \mathbb{R}^{6000 \times 5000}$. The response vector $\mathbf{y} \in \mathbb{R}^{6000 \times 1}$ consists of binary labels, indicating the class (either four or nine) of each sample.

Furthermore, we test screening algorithms on the sparse overlapping group lasso in (3.23) with $\lambda_1 = \lambda_2$ under different group structures. We chose this model because it induces individual level sparsity, allowing us to capture complex patterns of nonzero coefficients. Moreover, we generated group structures based on feature locality because features located nearby are often associated with an output vector jointly in image or genome datasets. The group structures used in our experiments are as follows:

1. $\ell_1 + \mathbf{nonoverlap}$ groups: \mathcal{G} contains nonoverlapping groups, where

³from <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

each group consists of 20 consecutive features.

2. $\ell_1 + \mathbf{tree\ structure\ groups}$: \mathcal{G} contains tree structured groups with four levels, where from the root to the leaves, the groups consist of 20/15/10/5 consecutive features. A parent group is always super set of their children groups.
3. $\ell_1 + \mathbf{overlap\ groups}$: \mathcal{G} contains overlapping groups, in which each group contains 20 features and consecutive groups overlap by 5 features.
4. $\ell_1 + \mathbf{overlap\ groups\ guided\ by\ prior\ knowledge}$: This group structure is constructed only for AD dataset. \mathcal{G} contains 26,222 groups, in which each group contains single nucleotide polymorphisms (i.e., features in AD dataset) located on the same gene region, defined by the interval between the gene’s transcription start and end site. Note that groups may overlap due to overlapping genes.

We solved the overlapping group lasso problems with OLS, SOLS, and GDPP on a sequence of $\{0.9^1\lambda', 0.9^2\lambda', \dots, 0.9^{30}\lambda'\}$, where λ' was obtained by Algorithm 3. In the sequential screening, we stop using screening from λ_{t+1} if no features are discarded at λ_t ($1 \leq t \leq 29$) because it is likely that screening with $\lambda < \lambda_t$ discards a few features, if at all. All the experiments except GISETTE (which contains a single \mathbf{y}) were repeated 10 times with a different \mathbf{y} for each run, and we report the average performance. As an overlapping group lasso solver, we used FoGLasso (Yuan, Liu and Ye, 2011), a state-of-the-art solver in the SLEP package (Liu, Ji and Ye, 2009), and all screening rules were implemented in Matlab. Below, we our demonstrate empirical results under different datasets, different group structures, and different group sizes to study how the screening efficiency changes under various scenarios.

4.1. *Different Datasets.* We first investigate the effectiveness of screening algorithms in terms of the screening rejection ratio under the four different datasets. For this experiment, we used the $\ell_1 + \mathbf{tree\ structure\ groups}$. Figures 2 (a–d) show the rejection ratio of OLS, SOLS, and GDPP on the four datasets as we change λ parameter. Figures 2 (e–h) demonstrate the differences between screening rejection ratios of the OLS and those of GDPP; each dot represents a result for a single λ , and the distance between dots and the diagonal line denotes the increased rejection ratio by OLS’s use of overlapping groups.

For all datasets except AD dataset, OLS and SOLS reject significantly more features than GDPP, taking advantage of overlapping groups to decrease the left-hand side of screening rules. However, given the AD dataset,

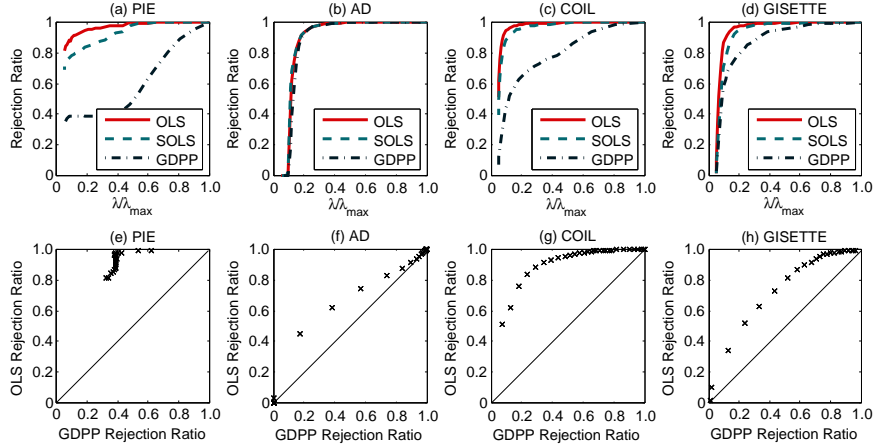


FIG 2. Rejection ratio of OLS (overlapping group lasso screening), SOLS (sparse overlapping group lasso screening) and GDPP (group dual polytope projection) and their comparison on (a,e) PIE, (b,f) COIL, (c,g) AD, and (d,h) GISETTE datasets for different λ/λ_{max} parameters.

the performance gap between OLS (or SOLS) and GDPP was the smallest. This phenomenon is due to the small groups used for the AD dataset. The median of AD group sizes was 8; in contrast, the other datasets used the groups, sizes of 20/15/10. In the experiments with different group sizes in §4.3, we confirm that the performance gap between OLS (or SOLS) and GDPP increases as the group size increases. This is not surprising because as group size increases, the number of free variables to minimize the left-hand side of screening rules increases, resulting in a better screening rejection ratio. Note that OLS and SOLS outperform GDPP for all $N \gg J$ (PIE), $N \ll J$ (COIL), and $N \approx J$ (GISETTE) settings. This observation suggests to us that OLS and SOLS would be useful for all large datasets regardless of their sample or feature sizes. Furthermore, OLS and SOLS show similar rejection ratios for all datasets. Therefore, in a single-machine setting, if ℓ_1 norm is included in the regularization, SOLS would be preferred over OLS due to SOLS’s low computational complexity as well as its rejection ratio comparable to that of OLS.

We then measured the speed gain achieved by the screening rules on the four different datasets. In Figure 3, we compare the running times of solving overlapping group lasso problems given a sequence of λ parameters among the solver without screening, the solver with OLS, the solver with SOLS, and the solver with OGDPP. For the results with screening, we also illustrate the portion of running times consumed by the solver and by screening. For

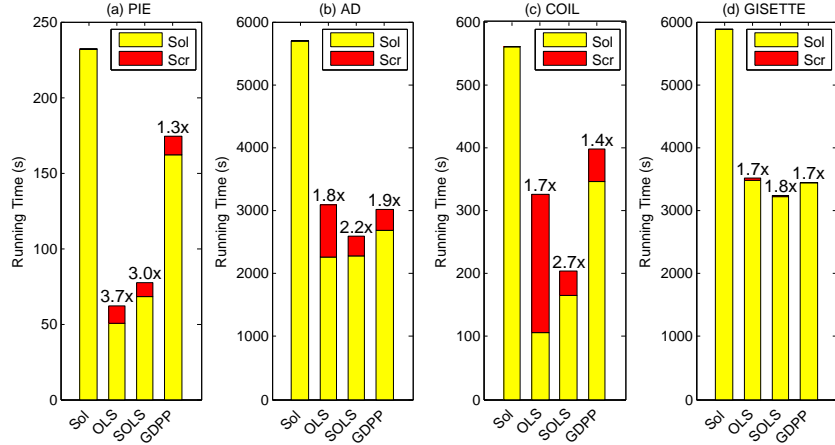


FIG 3. Running times of overlapping group lasso solver without screening (*Sol*), solver with OSL screening (*OSL*), solver with SOSL screening (*SOLS*), and solver with GDPP screening (*GDPP*) on (a) *PIE*, (b) *AD*, (c) *COIL*, and (d) *GISETTE* datasets.

the *PIE*, *AD*, *COIL*, and *GESETTE* datasets, we observed $3\times$, $2.2\times$, $2.7\times$, and $1.8\times$ speed-up by *SOLS* in comparison to the solver without screening. One interesting observation is that for $N \gg J$ setting such as in the *PIE* dataset, running times of screening with *OLS* and *SOLS* are similar; but for large J such as in the *COIL* dataset, *OLS* is significantly more expensive than *SOLS*. As a result, the solver with *OLS* was slower than the solver with *SOLS*, even though *OLS* discarded more features than *SOLS*. However, we note that screening rules are embarrassingly parallel. Thus, in a distributed setting, *OLS* would be appealing because the screening portion of running times can be reduced proportionally to the number of cores. Hereafter, we show only the screening rejection ratio by different screening algorithms because the similar patterns of running times are observed under the other experimental settings.

4.2. *Different Group Structures.* Next, we run the screening algorithms on the *COIL* dataset under three different group structures, including ℓ_1 + nonoverlap groups, ℓ_1 + tree structure groups, and ℓ_1 + overlap groups. Figures 4(a–c) show the rejection ratios of different algorithms as λ changes given different group structures; we compare the rejection ratio of *SOLS* with that of *GDPP* below each corresponding plot. Overall, we can see that *OLS* and *SOLS* maintain high rejection ratios over a wide range of λ parameters, but *GDPP*'s rejection ratio drops quickly as λ decreases. The rejection ratios by *OLS* and *SOLS* were identical under ℓ_1 + nonoverlap groups and ℓ_1 + overlap groups because only individual coefficients due to ℓ_1 penalty are

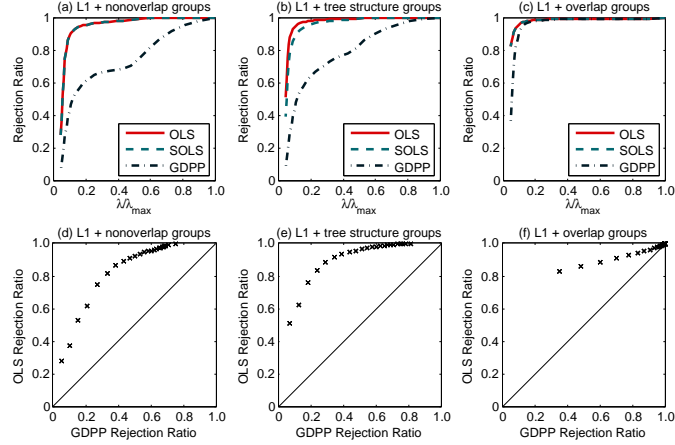


FIG 4. Rejection ratio of OLS, SOLS, and GDPP with different group structures: (a,d) ℓ_1 + nonoverlap groups, (b,e) ℓ_1 + tree structure groups, and (c,f) ℓ_1 + overlap groups.

inclusive groups for both screening methods. Interestingly, even under ℓ_1 + tree structure groups, OLS and SOLS show similar rejection ratios, which indicates that use of ℓ_1 penalty for $\tilde{\mathcal{G}}_1$ is effective.

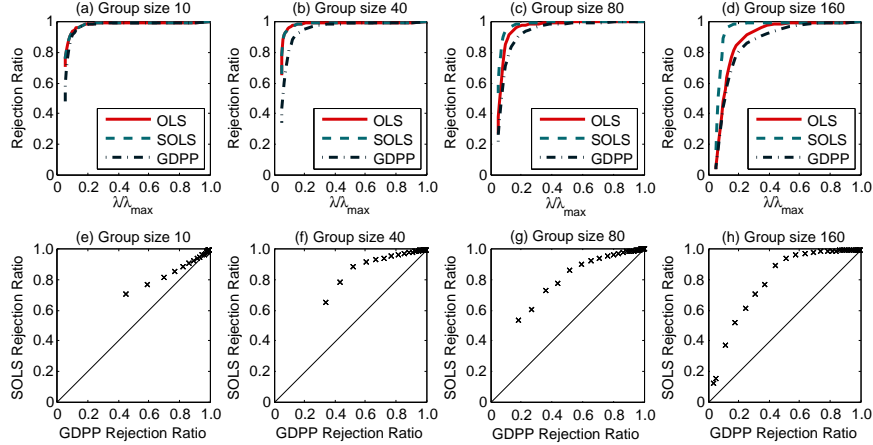


FIG 5. Rejection ratio of OGDPP and GPP with different sizes of overlapping groups. The sizes of groups are (a) 10, (b) 40, (c) 80, and (d) 160 features.

4.3. *Different Group Sizes.* We also observed the effects of group sizes on the screening rejection ratio. For this experiment, we used the COIL dataset and ℓ_1 + overlap groups, where the size of groups is changed from 10 to 160. Figure 5 shows the screening efficiency of GDPP and OGDPP

given different group sizes. Clearly, the rejection ratio of GDPP decreases as the group size increases. However, SOLS was not affected by the increased group sizes because SOLS can take advantage of the increased number of overlapping groups, $\bar{\mathcal{G}}_1 = \{j : j \in \mathbf{g}\}$, within each tested group \mathbf{g} . As the number of overlapping groups increases, the number of free variables to optimize in the left-hand side of SOLS rule also increases. Therefore, SOLS kept the high rejection ratios given all group sizes. The rejection ratio of OLS is decreased for group sizes ≥ 80 because we fixed the window size to search for $\bar{\mathcal{G}}_1$ by 50. Thus, OLS's rejection ratio started to decline from the group size of 80 due to the fixed number of overlapping groups in $\bar{\mathcal{G}}_1$.

5. Conclusion. In this paper, we developed screening rules including OLS and SOLS for overlapping group lasso. We make it possible to screen each group independently of each other by considering only groups inclusive of each tested group. Taking advantage of groups that overlap with tested groups, we showed that OLS and SOLS are efficient in terms of the screening rejection ratio. In addition, we verify that the GDPP screening rule (Wang et al., 2013) is a special case of OLS, but it is less efficient than OLS because it lacks the capability to use overlapping groups. In OLS, there is a step to find all groups overlapped with the group of interest, which can be computationally heavy. Thus, as a special case of OLS, we also present SOLS for sparse overlapping group lasso that includes ℓ_1 penalty. In our experiments, SOLS was much faster than OLS for screening, while maintaining its rejection ratio comparable to that of OLS. Motivated by enhanced DPP (Wang et al., 2013), which uses a tight range of dual optimal solutions, developing more efficient OLS or SOLS would be an interesting research direction. Furthermore, extending OLS to various loss functions such as logistic loss or hinge loss is left for future research.

APPENDIX A: DUAL FORMULATION OF OVERLAPPING GROUP LASSO

Overlapping group lasso problem is defined by

$$(A.1) \quad \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2,$$

where $\mathbf{X} \in \mathbb{R}^{N \times J}$ is the input data for J inputs and N samples, $\mathbf{y} \in \mathbb{R}^{N \times 1}$ is the output vector, $\boldsymbol{\beta} \in \mathbb{R}^{J \times 1}$ is the vector of regression coefficients, $n_{\mathbf{g}}$ is the size of group \mathbf{g} , and λ is a regularization parameter that determines the sparsity of $\boldsymbol{\beta}$. Here, we convert the primal overlapping group lasso problem in (A.1) to a dual problem.

Introducing $\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$, (A.1) can be written as

$$(A.2) \quad \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{z}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$$

subject to $\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$.

Lagrangian of (A.2) is

$$(A.3) \quad L(\boldsymbol{\beta}, \mathbf{z}, \boldsymbol{\eta}) = \frac{1}{2} \|\mathbf{z}\|_2^2 + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2 + \boldsymbol{\eta}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}).$$

Dual function $g(\boldsymbol{\eta})$ of (A.3) is

$$(A.4) \quad \begin{aligned} g(\boldsymbol{\eta}) &= \inf_{\boldsymbol{\beta}, \mathbf{z}} L(\boldsymbol{\beta}, \mathbf{z}, \boldsymbol{\eta}) \\ &= \boldsymbol{\eta}^T \mathbf{y} + \inf_{\boldsymbol{\beta}} \left(-\boldsymbol{\eta}^T \mathbf{X}\boldsymbol{\beta} + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2 \right) + \inf_{\mathbf{z}} \left(\frac{1}{2} \|\mathbf{z}\|_2^2 - \boldsymbol{\eta}^T \mathbf{z} \right). \end{aligned}$$

To obtain $g(\boldsymbol{\eta})$, we solve the following two optimization problems:

$$(A.5) \quad \inf_{\boldsymbol{\beta}} \left(-\boldsymbol{\eta}^T \mathbf{X}\boldsymbol{\beta} + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2 \right)$$

and

$$(A.6) \quad \inf_{\mathbf{z}} \left(\frac{1}{2} \|\mathbf{z}\|_2^2 - \boldsymbol{\eta}^T \mathbf{z} \right).$$

We first solve (A.5). Let us denote $f_1(\boldsymbol{\beta}) \equiv \left(-\boldsymbol{\eta}^T \mathbf{X}\boldsymbol{\beta} + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2 \right)$. A subgradient of $f_1(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$ is

$$(A.7) \quad \frac{\partial f_1(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -\mathbf{X}^T \boldsymbol{\eta} + \lambda \mathbf{v},$$

where \mathbf{v} is a subgradient of $\sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$ with respect to $\boldsymbol{\beta}$. The j -th element of \mathbf{v} is given by

$$(A.8) \quad v_j = \sum_{\{\mathbf{g}: j \in \mathbf{g}, \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \frac{\sqrt{n_{\mathbf{g}}} \beta_j}{\|\boldsymbol{\beta}_{\mathbf{g}}\|_2} + \sum_{\{\mathbf{g}: j \in \mathbf{g}, \boldsymbol{\beta}_{\mathbf{g}} = \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \sqrt{n_{\mathbf{g}}} o_j,$$

where o_j is a subgradient that satisfies $\|\mathbf{o}_{\mathbf{g}}\|_2 \leq 1$, where $j \in \mathbf{g}$. To obtain an optimal $\boldsymbol{\beta}^*$ for (A.5), we set $\frac{\partial f_1(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0$. Then, we have

$$(A.9) \quad \mathbf{X}^T \boldsymbol{\eta} = \lambda \mathbf{v}.$$

Plugging it into $f_1(\boldsymbol{\beta})$, we get

$$(A.10) \quad f_1(\boldsymbol{\beta}) = -\lambda \mathbf{v}^T \boldsymbol{\beta} + \lambda \sum_{\mathbf{g} \in \mathcal{G}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$$

$$(A.11) \quad = -\lambda \sum_j \sum_{\{\mathbf{g}: j \in \mathbf{g}, \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} v_j \beta_j + \lambda \sum_{\{\mathbf{g}: \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$$

$$(A.12) \quad = -\lambda \sum_j \sum_{\{\mathbf{g}: j \in \mathbf{g}, \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \frac{\sqrt{n_{\mathbf{g}}}(\beta_j)^2}{\|\boldsymbol{\beta}_{\mathbf{g}}\|_2} + \lambda \sum_{\{\mathbf{g}: \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$$

$$(A.13) \quad = -\lambda \sum_{\{\mathbf{g}: \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \sum_{j \in \mathbf{g}} \frac{\sqrt{n_{\mathbf{g}}}(\beta_j)^2}{\|\boldsymbol{\beta}_{\mathbf{g}}\|_2} + \lambda \sum_{\{\mathbf{g}: \boldsymbol{\beta}_{\mathbf{g}} \neq \mathbf{0}, \mathbf{g} \in \mathcal{G}\}} \sqrt{n_{\mathbf{g}}} \|\boldsymbol{\beta}_{\mathbf{g}}\|_2$$

$$(A.14) \quad = 0.$$

Here (3.10) is used for (A.12). Therefore, $\inf_{\boldsymbol{\beta}} f_1(\boldsymbol{\beta}) = 0$.

Now we solve the second problem $f_2(\mathbf{z}) \equiv \frac{1}{2} \|\mathbf{z}\|_2^2 - \boldsymbol{\eta}^T \mathbf{z}$. This result was presented in (Wang et al., 2013). However, here we show the derivation for self-containedness.

$$(A.15) \quad f_2(\mathbf{z}) = \frac{1}{2} \|\mathbf{z}\|_2^2 - \boldsymbol{\eta}^T \mathbf{z} = \frac{1}{2} \left(\|\mathbf{z} - \boldsymbol{\eta}\|_2^2 - \|\boldsymbol{\eta}\|_2^2 \right).$$

Note that $\boldsymbol{\eta} = \operatorname{argmin}_{\mathbf{z}} \frac{1}{2} \left(\|\mathbf{z} - \boldsymbol{\eta}\|_2^2 - \|\boldsymbol{\eta}\|_2^2 \right)$, and thus $\inf_{\mathbf{z}} f_2(\mathbf{z}) = -\frac{1}{2} \|\boldsymbol{\eta}\|_2^2$.

Based on these results for (A.5) and (A.6), the dual function $g(\boldsymbol{\eta})$ in (A.4) is given by,

$$(A.16) \quad g(\boldsymbol{\eta}) = \boldsymbol{\eta}^T \mathbf{y} - \frac{1}{2} \|\boldsymbol{\eta}\|_2^2 = \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{1}{2} \|\boldsymbol{\eta} - \mathbf{y}\|_2^2.$$

Finally, we denote $\boldsymbol{\theta} = \frac{\boldsymbol{\eta}}{\lambda}$. Combining (A.16) with (A.9), a dual formulation of overlapping group lasso is as follows:

$$(A.17) \quad \begin{aligned} & \sup_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{\lambda^2}{2} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2^2 \\ & \text{subject to } \mathbf{X}^T \boldsymbol{\theta} = \mathbf{v}. \end{aligned}$$

APPENDIX B: PROOF OF THEOREM 1

THEOREM 1. *For the overlapping lasso problem, suppose that we are given an optimal dual solution $\boldsymbol{\theta}^*(\lambda_0)$. Then for $\lambda < \lambda_0$, $\boldsymbol{\beta}_{\mathbf{g}}^*(\lambda) = \mathbf{0}$ if*

$$(B.1) \quad \min_{\mathbf{w}_{\mathbf{h}}, \|\mathbf{w}_{\mathbf{h}}\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \boldsymbol{\theta}^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_{\mathbf{h}}} w_j \right)^2} < \sqrt{n_{\mathbf{g}}} - \|\mathbf{X}_{\mathbf{g}}\|_F \|\mathbf{y}\|_2 \left| \frac{1}{\lambda} - \frac{1}{\lambda_0} \right|.$$

PROOF. Based on (3.14), we have a sphere Θ that contains $\theta^*(\lambda)$, which is centered at $\theta^*(\lambda_0)$ with a radius of $\rho = \left\| \frac{\mathbf{y}}{\lambda} - \frac{\mathbf{y}}{\lambda_0} \right\|_2$. Thus, we can represent $\theta^*(\lambda) = \theta^*(\lambda_0) + \mathbf{r}$, where $\|\mathbf{r}\|_2 \leq \rho$. Plugging it into (3.13) we get

$$\begin{aligned}
 \text{(B.2)} \quad b_{\mathbf{g}} &\leq \min_{\mathbf{w}_{\mathbf{h}}: \|\mathbf{w}_{\mathbf{h}}\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \theta^*(\lambda_0) + \mathbf{x}_j^T \mathbf{r} - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_{\mathbf{h}}} w_j \right)^2} \\
 \text{(B.3)} \quad &\leq \min_{\mathbf{w}_{\mathbf{h}}: \|\mathbf{w}_{\mathbf{h}}\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \theta^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_{\mathbf{h}}} w_j \right)^2} + \sqrt{\sum_{j \in \mathbf{g}} (\mathbf{x}_j^T \mathbf{r})^2} \\
 \text{(B.4)} \quad &\leq \min_{\mathbf{w}_{\mathbf{h}}: \|\mathbf{w}_{\mathbf{h}}\|_2 \leq 1} \sqrt{\sum_{j \in \mathbf{g}} \left(\mathbf{x}_j^T \theta^*(\lambda_0) - \sum_{j \in \mathbf{h}, \mathbf{h} \in \bar{\mathcal{G}}_1} \sqrt{n_{\mathbf{h}}} w_j \right)^2} + \sqrt{\|\mathbf{X}_{\mathbf{g}}\|_F^2 \|\mathbf{y}\|_2^2 \left(\frac{1}{\lambda} - \frac{1}{\lambda_0} \right)^2}.
 \end{aligned}$$

We used Minkowski's inequality for (B.3), and Cauchy-Schwarz inequality and $\|\mathbf{X}_{\mathbf{g}}\|_2^2 \leq \|\mathbf{X}_{\mathbf{g}}\|_F^2$ for (B.4). From (3.13), if $b_{\mathbf{g}} < \sqrt{n_{\mathbf{g}}}$, then $\beta_{\mathbf{g}}^* = \mathbf{0}$, and the result follows. \square

REFERENCES

- BACH, F., JENATTON, R., MAIRAL, J. and OBOZINSKI, G. (2012). Optimization with sparsity-inducing penalties. *Foundations and Trends[®] in Machine Learning* **4** 1–106.
- BERTSEKAS, D. P., NEDI, A., OZDAGLAR, A. E. et al. (2003). *Convex analysis and optimization*. Athena Scientific.
- BONNEFOY, A., EMIYA, V., RALAIVOLA, L., GRIBONVAL, R. et al. (2014). A Dynamic Screening test principle for the Lasso. In *European Signal Processing Conference*.
- CHEN, X., LIN, Q., KIM, S., CARBONELL, J. G., XING, E. P. et al. (2012). Smoothing proximal gradient method for general structured sparse regression. *The Annals of Applied Statistics* **6** 719–752.
- DENG, W., YIN, W. and ZHANG, Y. (2013). Group sparse optimization by alternating direction method. *Proceedings of the SPIE* **8858** 88580R.
- GHAOUI, L. E., VIALON, V. and RABBANI, T. (2012). Safe feature elimination in sparse supervised learning. *Pacific Journal of Optimization* **8** 667698.
- GUYON, I., GUNN, S., BEN-HUR, A. and DROR, G. (2004). Result analysis of the nips 2003 feature selection challenge. In *Advances in Neural Information Processing Systems* 545–552.
- JACOB, L., OBOZINSKI, G. and VERT, J. P. (2009). Group lasso with overlap and graph lasso. In *Proceedings of the 26th Annual International Conference on Machine Learning* 433–440. ACM.
- JENATTON, R., AUDIBERT, J.-Y. and BACH, F. (2011). Structured variable selection with sparsity-inducing norms. *The Journal of Machine Learning Research* **12** 2777–2824.
- KIM, S., XING, E. P. et al. (2012). Tree-guided group lasso for multi-response regression with structured sparsity, with an application to eQTL mapping. *The Annals of Applied Statistics* **6** 1095–1117.
- LECUN, Y., BOTTOU, L., BENGIO, Y. and HAFFNER, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86** 2278–2324.
- LEE, S. and XING, E. P. (2012). Leveraging input and output structures for joint mapping of epistatic and marginal eQTLs. *Bioinformatics* **28** i137–i146.

- LIU, J., JI, S. and YE, J. (2009). SLEP: Sparse learning with efficient projections. *Arizona State University* **6**.
- LIU, J., ZHAO, Z., WANG, J. and YE, J. (2013). Safe Screening With Variational Inequalities and Its Application to LASSO In *Proceedings of the 30th International Conference on Machine Learning*.
- NAYAR, S. K., NENE, S. A. and MURASE, H. (1996). Columbia object image library (coil 100). *Department of Comp. Science, Columbia University, Tech. Rep. CUCS-006-96*.
- SIM, T., BAKER, S. and BSAT, M. (2002). The CMU pose, illumination, and expression (PIE) database. In *Automatic Face and Gesture Recognition, 2002. Proceedings. Fifth IEEE International Conference on* 46–51. IEEE.
- SIMON, N., FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2013). A sparse-group lasso. *Journal of Computational and Graphical Statistics* **22** 231–245.
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* **58** 267–288.
- TIBSHIRANI, R., BIEN, J., FRIEDMAN, J., HASTIE, T., SIMON, N., TAYLOR, J. and TIBSHIRANI, R. J. (2012). Strong rules for discarding predictors in lasso-type problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **74** 245–266.
- WANG, J., ZHOU, J., WONKA, P. and YE, J. (2013). Lasso screening rules via dual polytope projection. In *Advances in Neural Information Processing Systems* 1070–1078.
- XIANG, Z. J. and RAMADGE, P. J. (2012). Fast lasso screening tests based on correlations. In *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on* 2137–2140. IEEE.
- XIANG, Z. J., XU, H. and RAMADGE, P. J. (2011). Learning sparse representations of high dimensional data on large scale dictionaries. In *Advances in Neural Information Processing Systems* 900–908.
- YANG, C., WAN, X., YANG, Q., XUE, H. and YU, W. (2010). Identifying main effects and epistatic interactions from large-scale SNP data via adaptive group Lasso. *BMC bioinformatics* **11** S18.
- YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *The Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **68** 49.
- YUAN, L., LIU, J. and YE, J. (2011). Efficient methods for overlapping group lasso. In *NIPS* 352–360.
- ZHANG, B., GAITERI, C., BODEA, L.-G., WANG, Z., MCELWEE, J., PODTELEZHNIKOV, A. A., ZHANG, C., XIE, T., TRAN, L., DOBRIN, R. et al. (2013). Integrated systems approach identifies genetic nodes and networks in late-onset Alzheimers disease. *Cell* **153** 707–720.
- ZHAO, P., ROCHA, G. and YU, B. (2009). The composite absolute penalties family for grouped and hierarchical variable selection. *The Annals of Statistics* **37** 3468–3497.