

Two-Layer Feature Reduction for Sparse-Group Lasso via Decomposition of Convex Sets

Authored by:

Jieping Ye
Jie Wang

Abstract

Sparse-Group Lasso (SGL) has been shown to be a powerful regression technique for simultaneously discovering group and within-group sparse patterns by using a combination of the l1 and l2 norms. However, in large-scale applications, the complexity of the regularizers entails great computational challenges. In this paper, we propose a novel two-layer feature reduction method (TLFre) for SGL via a decomposition of its dual feasible set. The two-layer reduction is able to quickly identify the inactive groups and the inactive features, respectively, which are guaranteed to be absent from the sparse representation and can be removed from the optimization. Existing feature reduction methods are only applicable for sparse models with one sparsity-inducing regularizer. To our best knowledge, TLFre is the first one that is capable of dealing with multiple sparsity-inducing regularizers. Moreover, TLFre has a very low computational cost and can be integrated with any existing solvers. Experiments on both synthetic and real data sets show that TLFre improves the efficiency of SGL by orders of magnitude.

1 Paper Body

Sparse-Group Lasso (SGL) [5, 16] is a powerful regression technique in identifying important groups and features simultaneously. To yield sparsity at both group and individual feature levels, SGL combines the Lasso [18] and group Lasso [28] penalties. In recent years, SGL has found great success in a wide range of applications, including but not limited to machine learning [20, 27], signal processing [17], bioinformatics [14] etc. Many research efforts have been devoted to developing efficient solvers for SGL [5, 16, 10, 21]. However, when the feature dimension is extremely high, the complexity of the SGL regularizers imposes great computational challenges. Therefore, there is an increasingly urgent need for nontraditional techniques to address the challenges posed by the massive volume of the data sources. Recently, El Ghaoui et al. [4] proposed a promising feature reduction method, called SAFE screening, to screen out

the so-called inactive features, which have zero coefficients in the solution, from the optimization. Thus, the size of the data matrix needed for the training phase can be significantly reduced, which may lead to substantial improvement in the efficiency of solving sparse models. Inspired by SAFE, various exact and heuristic feature screening methods have been proposed for many sparse models such as Lasso [25, 11, 19, 26], group Lasso [25, 22, 19], etc. It is worthwhile to mention that the discarded features by exact feature screening methods such as SAFE [4], DOME [26] and EDPP [25] are guaranteed to have zero coefficients in the solution. However, heuristic feature screening methods like Strong Rule [19] may mistakenly discard features which have nonzero coefficients in the solution. More recently, the idea of exact feature screening has been extended to exact sample screening, which screens out the nonsupport vectors in SVM [13, 23] and LAD [23]. As a promising data reduction tool, exact feature/sample screening would be of great practical importance because they can effectively reduce the data size without sacrificing the optimality [12]. 1

However, all of the existing feature/sample screening methods are only applicable for the sparse models with one sparsity-inducing regularizer. In this paper, we propose an exact two-layer feature screening method, called TLFre, for the SGL problem. The two-layer reduction is able to quickly identify the inactive groups and the inactive features, respectively, which are guaranteed to have zero coefficients in the solution. To the best of our knowledge, TLFre is the first screening method which is capable of dealing with multiple sparsity-inducing regularizers. We note that most of the existing exact feature screening methods involve an estimation of the dual optimal solution. The difficulty in developing screening methods for sparse models with multiple sparsity-inducing regularizers like SGL is that the dual feasible set is the sum of simple convex sets. Thus, to determine the feasibility of a given point, we need to know if it is decomposable with respect to the summands, which is itself a nontrivial problem (see Section 2). One of our major contributions is that we derive an elegant decomposition method of any dual feasible solutions of SGL via the framework of Fenchel's duality (see Section 3). Based on the Fenchel's dual problem of SGL, we motivate TLFre by an in-depth exploration of its geometric properties and the optimality conditions. We derive the set of the regularization parameter values corresponding to zero solutions. To develop TLFre, we need to estimate the upper bounds involving the dual optimal solution. To this end, we first give an accurate estimation of the dual optimal solution via the normal cones. Then, we formulate the estimation of the upper bounds via nonconvex optimization problems. We show that these nonconvex problems admit closed form solutions. Experiments on both synthetic and real data sets demonstrate that the speedup gained by TLFre in solving SGL can be orders of magnitude. All proofs are provided in the long version of this paper [24].

Notation: Let $\|\cdot\|_1$, $\|\cdot\|_2$ and $\|\cdot\|_\infty$ be the '1', '2' and ' ∞ ' norms, respectively. Denote by B_1^n , B_2^n , and B_∞^n the unit '1', '2', and ' ∞ ' norm balls in \mathbb{R}^n (we omit the superscript if it is clear from the context). For a set C , let $\text{int } C$ be its interior. If C is closed and convex, we define the projection operator as $\text{PC}(w) := \arg\min_{C} \|w - u\|_2$. We denote by $\text{IC}(C)$ the indicator function of C , which is 0 on C and

elsewhere. Let $\mathcal{F}_0(\mathbb{R}^n)$ be the class of proper closed convex functions on \mathbb{R}^n . For $f \in \mathcal{F}_0(\mathbb{R}^n)$, let ∂f be its subdifferential. The domain of f is the set $\text{dom } f := \{w : f(w) < \infty\}$. For $w \in \mathbb{R}^n$, let $[w]_i$ be its i th component. For $\gamma \in \mathbb{R}$, let $\text{sgn}(\gamma) = \text{sign}(\gamma)$ if $\gamma \neq 0$, and $\text{sgn}(0) = 0$. We define $\text{sign}([w]_i)$, if $[w]_i \neq 0$; $\text{SGN}(w) = s \in \mathbb{R}^n : [s]_i \in [-1, 1]$, if $[w]_i = 0$. We denote by $\gamma_+ = \max(\gamma, 0)$. Then, the shrinkage operator $S_\gamma(w) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ with $\gamma \geq 0$ is $[S_\gamma(w)]_i = (-[w]_i - \gamma)_+ \text{sgn}([w]_i)$, $i = 1, \dots, n$.

(1)

Basics and Motivation

In this section, we briefly review some basics of SGL. Let $y \in \mathbb{R}^n$ be the response vector and $X \in \mathbb{R}^n \times \mathbb{R}^p$ be the matrix of features. With the group information available, the SGL problem [5] is

$$\min_{\beta} \|X\beta - y\|_2^2 + \sum_{g=1}^G \lambda_g \|\beta_g\|_1, \quad (2)$$

where

$$\beta_g = (X_g \beta_g)_{g=1}^G, \quad \lambda_g \geq 0, \quad (2)$$

where λ_g is the number of features in the g th group, $X_g \in \mathbb{R}^{n_g \times p_g}$ denotes the predictors in that group with the corresponding coefficient vector β_g , and λ_1, λ_2 are positive regularization parameters. Without loss of generality, let $\lambda_1 = 1$ and $\lambda_2 = \lambda$ with $\lambda \geq 0$. Then, problem (2) becomes:

$$\min_{\beta} \|X\beta - y\|_2^2 + \|\beta\|_1 + \lambda \|\beta\|_2, \quad (3)$$

By the Lagrangian multipliers method [24], the dual problem of SGL is

$$\max_{\gamma} \sum_{g=1}^G \lambda_g \|\gamma_g\|_2^2 - \frac{1}{2} \|X^T \gamma - y\|_2^2, \quad \gamma_g \in \mathbb{R}^{p_g}, \quad g = 1, \dots, G. \quad (4)$$

It is well-known that the dual feasible set of Lasso is the intersection of closed half spaces (thus a polytope); for group Lasso, the dual feasible set is the intersection of ellipsoids. Surprisingly, the geometric properties of these dual feasible sets play fundamentally important roles in most of the existing screening methods for sparse models with one sparsity-inducing regularizer [23, 11, 25, 4]. When we incorporate multiple sparse-inducing regularizers to the sparse models, problem (4) indicates that the dual feasible set can be much more complicated. Although (4) provides a geometric

description of the dual feasible set of SGL, it is not suitable for further analysis. Notice that, even the feasibility of a given point γ is not easy to determine, since it is nontrivial to tell if $X^T \gamma - y$ can be decomposed into $b_1 + b_2$ with $b_1 \in \sum_{g=1}^G B_g$ and $b_2 \in B$. Therefore, to develop screening methods for SGL, it is desirable to gain deeper understanding of the sum of simple convex sets. In the next section, we analyze the dual feasible set of SGL in depth via the Fenchel's Duality Theorem. We show that for each $X^T \gamma - y \in D_g$, Fenchel's duality naturally leads to an explicit decomposition $X^T \gamma - y = b_1 + b_2$, with

one belonging to $\text{int } B$ and the other one belonging to $B^?$. This lays the foundation of the proposed screening method for SGL.

3

The Fenchel's Dual Problem of SGL

In Section 3.1, we derive the Fenchel's dual of SGL via Fenchel's Duality Theorem. We then motivate TLFre and sketch our approach in Section 3.2. In Section 3.3, we discuss the geometric properties of the Fenchel's dual of SGL and derive the set of (λ, μ) leading to zero solutions.

3.1 The Fenchel's Dual of SGL via Fenchel's Duality Theorem

To derive the Fenchel's dual problem of SGL, we need the Fenchel's Duality Theorem as stated in Theorem 1. The conjugate of $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is the function $f^* : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by $f^*(z) = \sup_{w \in \mathbb{R}^n} \langle z, w \rangle - f(w)$. Theorem 1. [Fenchel's Duality Theorem] Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$, and $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an affine mapping from \mathbb{R}^n to \mathbb{R}^m . Let p^* , d^* be primal and dual values defined, respectively, by the Fenchel problems: $p^* = \inf_{x \in \mathbb{R}^n} f(x) + g(Tx)$; $d^* = \sup_{\lambda \in \mathbb{R}^m} \lambda^T (b - \sum_{i=1}^n \lambda_i A_i x_i)$. One has $p^* \leq d^*$. If, furthermore, f and g satisfy the condition $0 \in \text{int}(\text{dom } f + T \text{dom } g)$, then the equality holds, i.e., $p^* = d^*$, and the supreme is attained in the dual problem if finite. We omit the proof of Theorem 1 since it is a slight modification of Theorem 3.3.5 in [2]. Let $f(w) = \frac{1}{2} \|w\|_2^2$, and $g(y) = \sum_{k=1}^K \max\{0, y_k\}$ be the second term in (3). Then, SGL can be written as $\min_x f(y - Xx) + g(y)$. To derive the Fenchel's dual problem of SGL, Theorem 1 implies that we need to find f^* and g^* . It is well-known that $f^*(z) = \frac{1}{2} \|z\|_2^2$. Therefore, we only need to find g^* , where the concept infimal convolution is needed. Let $h, g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$. The infimal convolution of h and g is defined by $(hg)(?) = \inf_{z \in \mathbb{R}^n} h(z) + g(? - z)$, and it is exact at a point $?$ if there exists a $z \in \mathbb{R}^n$ such that $(hg)(?) = h(z) + g(? - z)$. hg is exact if it is exact at every point of its domain, in which case it is denoted by $h \circ g$. Lemma 2. Let $g_1(?) = \sum_{k=1}^K \max\{0, y_k\}$, $g_2(?) = \sum_{k=1}^K \max\{0, y_k\}$ and $g(?) = g_1(?) + g_2(?)$. Moreover, let $Cg = \sum_{g=1}^G B_g^T Rng$, $g = 1, \dots, G$. Then, the following hold: (i): $(g_1 \circ g_2)(?) = \sum_{g=1}^G \max\{0, y_g\}$, $g_1 \circ g_2 = \sum_{g=1}^G \max\{0, y_g\}$,

(ii): $(g_1 \circ g_2)(?) = (\sum_{g=1}^G \max\{0, y_g\}) (?) = \sum_{g=1}^G \max\{0, y_g\}$ where $g \in Rng$ is the sub-vector of y corresponding to the g th group. Note that $PB^T (y_g)$ admits a closed form solution, i.e., $[PB^T (y_g)]_i = \text{sgn}([y_g]_i) \min(|[y_g]_i|, 1)$. Combining Theorem 1 and Lemma 2, the Fenchel's dual of SGL can be derived as follows. Theorem 3. For the SGL problem in (3), the following hold: (i): The Fenchel's dual of SGL is given by:

$$\max_{\lambda} \sum_{k=1}^K \lambda_k y_k - \sum_{k=1}^K \lambda_k^2 : \sum_{k=1}^K \lambda_k = \sum_{k=1}^K \lambda_k^2, \lambda_k \geq 0, g = 1, \dots, G. \quad (5)$$

(ii): Let $\lambda^*(?, ?)$ and $\lambda^{**}(?, ?)$ be the optimal solutions of problems (3) and (5), respectively. Then, $\lambda^{**}(?, ?) = y - X\lambda^*(?, ?)$, $\sum_{k=1}^K \lambda_k^* = \sum_{k=1}^K \lambda_k^{**}$, $g = 1, \dots, G$.

(6) (7)

Remark 1. We note that the shrinkage operator can also be expressed by $S^g(w) = w - P^g B^g(w)$, $g \geq 0$.

(8)

Therefore, problem (5) can be written more compactly as

$$\inf_{\mathbf{y}} \sum_{k=1}^K \sum_{g=1}^G \mathbf{1}_{\{\mathbf{y}^T \mathbf{X} \mathbf{T}_g \leq -\mathbf{b}_g\}} \quad \text{s.t. } \mathbf{y} \in \mathcal{G}.$$

(9)

?

Remark 2. Eq. (6) and Eq. (7) can be obtained by the Fenchel-Young inequality [2, 24]. They are the so-called KKT conditions [3] and can also be obtained by the Lagrangian multiplier method [24]. Moreover, for the SGL problem, its Lagrangian dual in (4) and Fenchel's dual in (5) are indeed equivalent to each other [24]. Remark 3. An appealing advantage of the Fenchel's dual in (5) is that we have a natural decomposition of all points $\mathbf{y}^T \mathbf{D}_g : \mathbf{y}^T \mathbf{D}_g = \mathbf{y}^T \mathbf{P}_g + \mathbf{y}^T \mathbf{S}_1(\mathbf{y})$ with $\mathbf{y}^T \mathbf{P}_g \leq \mathbf{b}_g$ and $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$. As a result, this leads to a convenient way to determine the feasibility of any dual variable \mathbf{y} by checking if $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$, $g = 1, \dots, G$. 3.2 Motivation of the Two-Layer Screening Rules We motivate the two-layer screening rules via the KKT condition in Eq. (7). As implied by the name, there are two layers in our method. The first layer aims to identify the inactive groups, and the second layer is designed to detect the inactive features for the remaining groups. by Eq. (7), we have the following cases by noting $\mathbf{1}_{\{\mathbf{y}^T \mathbf{D}_g \leq -\mathbf{b}_g\}} = \text{SGN}(\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g)$ and $\mathbf{0} \leq \mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g \leq 1$, if $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g = 0$, $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g = \{u : u \in [0, 1]\}$, if $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g = 0$. Case 1. If $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g = 0$, we have $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$ if and only if $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$, if $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g = 0$, $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$ if and only if $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$.

(10)

In view of Eq. (10), we can see that $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g \leq 0$ (a): $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$ if and only if $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$ and $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$.

(b): If $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g \leq 0$ then $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$. Case 2. If $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g > 0$

(11) (12)

$= 0$, we have $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$ if and only if $\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g$.

(13)

The first layer (group-level) of TLFr From (11) in Case 1, we have

$$\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g \quad \text{s.t. } \mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g \leq 0.$$

(R1)

Clearly, (R1) can be used to identify the inactive groups and thus a group-level screening rule. The second layer (feature-level) of TLFr Let x_{gi} be the i th column of \mathbf{X}_g . $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g = \mathbf{x}_{gi}^T \mathbf{y}$. In view of (12) and (13), we can see that $\mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g \leq 0$ if and only if $\mathbf{x}_{gi}^T \mathbf{y} \leq -\mathbf{b}_g$.

We have (R2)

Different from (R1), (R2) detects the inactive features and thus it is a feature-level screening rule. However, we cannot directly apply (R1) and (R2) to identify the inactive groups/features because both need to know \mathbf{y} . Inspired by the SAFE rules [4], we can first estimate a region \mathcal{R} containing \mathbf{y} . Let $\mathbf{X}_g = \{\mathbf{X}_g : \mathbf{y} \in \mathcal{R}\}$. Then, (R1) and (R2) can be relaxed as follows:

$$\mathbf{y}^T \mathbf{S}_1(\mathbf{y}) \leq \mathbf{c}_g \quad \text{s.t. } \mathbf{y}^T \mathbf{D}_g + \mathbf{b}_g \leq 0, \quad \mathbf{y} \in \mathcal{R}.$$

$$\sup_{\mathbf{y}} \mathbf{x}^T \mathbf{g} : \mathbf{y} \in \mathcal{F}(\mathbf{g}) \text{ if } \mathbf{g}^T \mathbf{y} = 0. \quad (\text{R2}^*)$$

Inspired by (R1^{*}) and (R2^{*}), we develop TLFRe via the following three steps: Step 1. Given \mathbf{g} and \mathbf{y} , we estimate a region $\mathcal{F}(\mathbf{g})$ that contains $\mathcal{F}(\mathbf{y})$. Step 2. We solve for the supreme values in (R1^{*}) and (R2^{*}). Step 3. By plugging in the supreme values from Step 2, (R1^{*}) and (R2^{*}) result in the desired two-layer screening rules for SGL. 4

3.3

The Set of Parameter Values Leading to Zero Solution ? For notational convenience, let $\mathcal{F}_g = \{\mathbf{y} : \mathbf{y}^T \mathbf{g} \leq \mathbf{y}^T \mathbf{y}\}$, $g = 1, \dots, G$; and thus the feasible set of the Fenchel's dual of SGL is $\mathcal{F} = \bigcap_{g=1, \dots, G} \mathcal{F}_g$. In view of problem (5) [or (9)], we can see that $\mathcal{F}(\mathbf{y})$ is the projection of $\mathbf{y}/\|\mathbf{y}\|$ on \mathcal{F} , i.e., $\mathcal{F}(\mathbf{y}) = \text{PF}(\mathbf{y}/\|\mathbf{y}\|)$. Thus, if $\mathbf{y}/\|\mathbf{y}\| \in \mathcal{F}$, we have $\mathcal{F}(\mathbf{y}) = \mathbf{y}/\|\mathbf{y}\|$. Moreover, by (R1), we can see that $\mathbf{y}^T \mathbf{g} = 0$ if $\mathbf{y}/\|\mathbf{y}\|$ is an interior point of \mathcal{F} . Indeed, we have the following stronger result.

THEOREM 4

Theorem 4. For the SGL problem, let $\mathbf{y}^* = \max_{\mathbf{y} \in \mathcal{F}} \{\mathbf{y}^T \mathbf{g} : \mathbf{y} \in \mathcal{F}_g\}$. Then, $\mathbf{y}^* \in \mathcal{F}$ if and only if $\mathbf{y}^* \in \mathcal{F}$. \mathbf{y}^* in the definition of \mathbf{y}^* admits a closed form solution [24]. Theorem 4 implies that the optimal solution \mathbf{y}^* is 0 as long as $\mathbf{y}/\|\mathbf{y}\| \in \mathcal{F}$. This geometric property also leads to an explicit characterization of the set of $(\mathbf{y}_1, \mathbf{y}_2)$ such that the corresponding solution of problem (2) is 0. We denote by $\mathcal{Z}(\mathbf{y}_1, \mathbf{y}_2)$ the optimal solution of problem (2). Corollary 5. For the SGL problem in (2), let $\mathbf{y}^* = \max_{\mathbf{y} \in \mathcal{F}} \{\mathbf{y}^T \mathbf{g} : \mathbf{y} \in \mathcal{F}_g\}$. Then, 1 (i): $\mathcal{Z}(\mathbf{y}_1, \mathbf{y}_2) = 0$ if $\mathbf{y}_1 \in \mathcal{F}$. 1

(ii): If $\mathbf{y}_1 \in \mathcal{F}$, $\mathbf{y}_2 \in \mathcal{F}$, then $\mathbf{y}^* = \max_{\mathbf{y} \in \mathcal{F}} \{\mathbf{y}^T \mathbf{g} : \mathbf{y} \in \mathcal{F}_g\}$.

4

$\mathbf{y}_1^T \mathbf{g} \leq \mathbf{y}_2^T \mathbf{g}$

or $\mathbf{y}_2^T \mathbf{g} \leq \mathbf{y}_1^T \mathbf{g}$, then $\mathcal{Z}(\mathbf{y}_1, \mathbf{y}_2) = 0$. 2

The Two-Layer Screening Rules for SGL

We follow the three steps in Section 3.2 to develop TLFRe. In Section 4.1, we give an accurate estimation of $\mathcal{F}(\mathbf{y})$ via normal cones [15]. Then, we compute the supreme values in (R1^{*}) and (R2^{*}) by solving nonconvex problems in Section 4.2. We present the TLFRe rules in Section 4.3. 4.1 Estimation of the Dual Optimal Solution Because of the geometric property of the dual problem in (5), i.e., $\mathcal{F}(\mathbf{y}) = \text{PF}(\mathbf{y}/\|\mathbf{y}\|)$, we have a very useful characterization of the dual optimal solution via the so-called normal cones [15]. Definition 1. [15] For a closed convex set $C \subset \mathbb{R}^n$ and a point $\mathbf{w} \in C$, the normal cone to C at \mathbf{w} is $\text{NC}(\mathbf{w}) = \{\mathbf{v} : \mathbf{v}^T \mathbf{w} \geq \mathbf{v}^T \mathbf{w}_i, \forall \mathbf{w}_i \in C\}$. (14) $\mathbf{v} \in \text{NC}(\mathbf{w})$ is known if $\mathbf{v}^T \mathbf{w} = \max_{\mathbf{w}_i \in C} \mathbf{v}^T \mathbf{w}_i$. Thus, we can estimate $\mathcal{F}(\mathbf{y})$ in terms of $\mathcal{F}(\mathbf{y})$. By Theorem 4, $\mathcal{F}(\mathbf{y}) = \mathcal{F}$ for $\mathbf{y}/\|\mathbf{y}\| \in \mathcal{F}$ to be estimated. Due to the same reason, we only consider the cases with $\mathbf{y}/\|\mathbf{y}\| \notin \mathcal{F}$. Remark 4. In many applications, the parameter values that perform the best are usually unknown. To determine appropriate parameter values, commonly used approaches such as cross validation and stability selection involve solving SGL many times over a grip of parameter values. Thus, given $\{\mathbf{y}_i\}_{i=1}^J$ and $\mathcal{Z}(\mathbf{y}_1, \dots, \mathbf{y}_J)$, we can fix the value of \mathbf{y} each time and solve SGL by varying the value of \mathbf{g} .

0.1 0.2 0.4 ?/??max
 1
 0.7 0.5 0.3 0.1 0.1 0.2 0.4 ?/??max
 1
 0.7 0.5 0.3
 0.1 0.2 0.4 ?/??max
 1
 ?
 (e) ? = $\tan(45^\circ)$
 1
 0.5 0.3
 0.01 0.02 0.04
 0.7 0.5 0.3
 0.1 0.2 0.4 ?/??max
 1
 1 0.9 0.7 0.5 0.3 0.1
 0.1 0.2 0.4 ?/??max
 1
 0.01 0.02 0.04
 ?
 (f) ? = $\tan(60^\circ)$
 0.7
 (d) ? = $\tan(30^\circ)$
 1 0.9
 0.01 0.02 0.04
 1 0.9
 0.1 0.1 0.2 0.4 ?/??max
 0.1
 0.01 0.02 0.04
 ?
 0.3
 (c) ? = $\tan(15^\circ)$
 1 0.9
 0.1
 0.01 0.02 0.04
 0.5
 0.01 0.02 0.04
 (b) ? = $\tan(5^\circ)$
 (a) 1 0.9
 0.7
 0.1
 Rejection Ratio
 0 0
 0.5
 Rejection Ratio
 100

0.7
 1 0.9
 Rejection Ratio
 200
 1 0.9
 Rejection Ratio
 ?1
 300
 Rejection Ratio
 ?max 1 (?2) ? = tan(5?) ? = tan(15?) ? = tan(30?) ? = tan(45?) ?
 = tan(60?) ? = tan(75?) ? = tan(85?)
 400
 0.1 0.2 0.4 ?/?max
 1
 ?
 (g) ? = tan(75)
 (h) ? = tan(85)
 Figure 1: Rejection ratios of TLFre on the Synthetic 1 data set.
 0.5 0.3 0.1
 500 ?2
 0.01 0.02 0.04
 1000
 Rejection Ratio
 Rejection Ratio
 0.7 0.5 0.3 0.1 0.01 0.02 0.04
 0.1 0.2 0.4 ?/?max
 1
 ?
 (e) ? = tan(45)
 1
 0.3
 0.7 0.5 0.3
 1
 0.1 0.2 0.4 ?/?max
 1
 0.7 0.5 0.3
 0.01 0.02 0.04
 ?
 0.5 0.3
 0.01 0.02 0.04
 0.1 0.2 0.4 ?/?max
 1
 1 0.9 0.7 0.5 0.3 0.1
 0.1 0.2 0.4 ?/?max
 ?
 (f) ? = tan(60)

0.7
 (d) $\theta = \tan(30^\circ)$
 1 0.9
 0.1
 0.01 0.02 0.04
 1 0.9
 0.1 0.1 0.2 0.4 θ/θ_{\max}
 (c) $\theta = \tan(15^\circ)$
 1 0.9
 0.1 0.1 0.2 0.4 θ/θ_{\max}
 0.5
 0.01 0.02 0.04
 (b) $\theta = \tan(5^\circ)$
 (a) 1 0.9
 0.7
 0.1
 Rejection Ratio
 0 0
 Rejection Ratio
 0.7
 1 0.9
 Rejection Ratio
 200
 1 0.9
 Rejection Ratio
 θ_1
 400
 Rejection Ratio
 $\theta_{\max} 1 (2^\circ) \theta = \tan(5^\circ) \theta = \tan(15^\circ) \theta = \tan(30^\circ) \theta = \tan(45^\circ) \theta$
 $= \tan(60^\circ) \theta = \tan(75^\circ) \theta = \tan(85^\circ)$
 600
 (g) $\theta = \tan(75^\circ)$
 1
 0.01 0.02 0.04
 0.1 0.2 0.4 θ/θ_{\max}
 1
 ?
 (h) $\theta = \tan(85^\circ)$

Figure 2: Rejection ratios of TLFre on the Synthetic 2 data set. 5.1 Simulation Studies We perform experiments on two synthetic data sets that are commonly used in the literature [19, 29]. The true model is $y = X\theta + 0.01$, $\theta \sim N(0, 1)$. We generate two data sets with 250 \times 10000 entries: Synthetic 1 and Synthetic 2. We randomly break the 10000 features into 1000 groups. For Synthetic 1, the entries of the data matrix X are i.i.d. standard Gaussian with pairwise correlation zero, i.e., $\text{corr}(x_i, x_j) = 0$. For Synthetic 2, the entries of the data matrix X are drawn from i.i.d. standard Gaussian with pairwise

correlation $0.5 - i/j$, i.e., $\text{corr}(x_i, x_j) = 0.5 - i/j$. To construct \mathbf{X} , we first randomly select γ_1 percent of groups. Then, for each selected group, we randomly select γ_2 percent of features. The selected components of \mathbf{X} are populated from a standard Gaussian and the remaining ones are set to 0. We set $\gamma_1 = \gamma_2 = 10$ for Synthetic 1 and $\gamma_1 = \gamma_2 = 20$ for Synthetic 2. The figures in the upper left corner of Fig. 1 and Fig. 2 show the plots of $\gamma_{\max}(\gamma_2)$ (see Corollary 1.5) and the sampled parameter values of γ_1 and γ_2 (recall that $\gamma_1 = \gamma_2$ and $\gamma_2 = \gamma$). For the other figures, the blue and red regions represent the rejection ratios of (L1) and (L2), respectively. We can see that TLFre is very effective in discarding inactive groups/features; that is, more than 90% of inactive features can be detected. Moreover, we can observe that the first layer screening (L1) becomes more effective with a larger γ . Intuitively, this is because the group Lasso penalty plays a more important role in enforcing the sparsity with a larger value of γ (recall that $\gamma_1 = \gamma_2$). The top and middle parts of Table 1 indicate that the speedup gained by TLFre is very significant (up to 30 times) and TLFre is very efficient. Compared to the running time of the solver without screening, the running time of TLFre is negligible. The running time of TLFre includes that of computing $\mathbf{X}^T \mathbf{X} \mathbf{g}$, $\mathbf{g} = 1, \dots, G$, which can be efficiently computed by the power method [6]. Indeed, this can be shared for TLFre with different parameter values.

5.2 Experiments on Real Data Set

We perform experiments on the Alzheimer’s Disease Neuroimaging Initiative (ADNI) data set (<http://adni.loni.usc.edu/>). The data matrix consists of 747 samples with 426040 single

Table 1: Running time (in seconds) for solving SGL along a sequence of 100 tuning parameter values of γ equally spaced on the logarithmic scale of γ/γ_{\max} max from 1.0 to 0.01 by (a): the solver [9] without screening; (b): the solver combined with TLFre. The top and middle parts report the results of TLFre on Synthetic 1 and Synthetic 2. The bottom part reports the results of TLFre on the ADNI data set with the GMV data as response.

γ	solver	TLFre Synthetic 1	TLFre+solver speedup
298.36	0.77	10.26	29.09
301.74	0.78	12.47	24.19
308.69	0.79	15.73	19.63
307.71	0.79	17.69	17.40
311.33	0.81	19.71	15.79
307.53	0.79	21.95	14.01
291.24	0.77	22.53	12.93
γ	solver	TLFre Synthetic 2	TLFre+solver speedup
294.64	0.79	11.05	26.66
294.92	0.80	12.89	22.88
297.29	0.80	16.08	18.49
297.50	0.81	18.90	15.74
297.59	0.81	20.45	14.55
295.51	0.81	21.58	13.69
292.24	0.82	22.80	12.82

30838.29 64.96 386.80 79.73
 31096.10 65.00 402.72 77.22
 30850.78 64.89 391.63 78.78
 30728.27 65.17 385.98 79.61
 30572.35 65.05 382.62 79.90
 0.5 0.3 0.1
 50
 100 ?2
 150
 1
 0.7 0.5 0.3 0.1 0.1 0.2 0.4 ?/??max
 (e) ? = tan(45?)
 1
 0.3
 1
 0.7 0.5 0.3
 0.1 0.2 0.4 ?/??max
 1
 (f) ? = tan(60?)
 0.3
 0.1 0.2 0.4 ?/??max
 1
 (d) ? = tan(30)
 1 0.9 0.7 0.5 0.3
 0.01 0.02 0.04
 0.5
 ?
 0.1
 0.01 0.02 0.04
 0.7
 0.01 0.02 0.04
 (c) ? = tan(15)
 1 0.9
 1 0.9
 0.1 0.1 0.2 0.4 ?/??max
 ?
 0.1
 0.01 0.02 0.04
 0.5
 0.01 0.02 0.04
 (b) ? = tan(5) Rejection Ratio
 Rejection Ratio
 0.1 0.2 0.4 ?/??max
 ?
 (a)
 0.7

0.1
 0.01 0.02 0.04
 1 0.9
 1 0.9
 Rejection Ratio
 0.7
 Rejection Ratio
 0 0
 1 0.9
 Rejection Ratio
 50
 Rejection Ratio
 $\theta_{\max} 1 (\theta_2) \theta = \tan(5^\circ) \theta = \tan(15^\circ) \theta = \tan(30^\circ) \theta = \tan(45^\circ) \theta$
 $= \tan(60^\circ) \theta = \tan(75^\circ) \theta = \tan(85^\circ) \theta$
 100
 θ_1
 solver 30652.56 30755.63 TLFre 64.08 64.56 TLFre+solver 372.04 383.17
 speedup 82.39 80.27
 Rejection Ratio
 ADNI+GMV
 1 0.9 0.7 0.5 0.3 0.1
 0.1 0.2 0.4 θ/θ_{\max}
 $(g) \theta = \tan(75^\circ)$
 1
 0.01 0.02 0.04
 0.1 0.2 0.4 θ/θ_{\max}
 1
 $(h) \theta = \tan(85^\circ)$

Figure 3: Rejection ratios of TLFre on the ADNI data set with grey matter volume as response. nucleotide polymorphisms (SNPs), which are divided into 94765 groups. The response vector is the grey matter volume (GMV). The figure in the upper left corner of Fig. 3 shows the plots of $\theta_{\max} (\theta_2)$ (see Corollary 5) and the 1 sampled parameter values of θ and θ_1 . The other figures present the rejection ratios of (L1) and (L2) by blue and red regions, respectively. We can see that almost all of the inactive groups/features are discarded by TLFre. The rejection ratios of $r_1 + r_2$ are very close to 1 in all cases. The bottom part of Table 1 shows that TLFre leads to a very significant speedup (about 80 times). In other words, the solver without screening needs about eight and a half hours to solve the 100 SGL problems for each value of θ . However, combined with TLFre, the solver needs only six to eight minutes. Moreover, we can observe that the computational cost of TLFre is negligible compared to that of the solver without screening. This demonstrates the efficiency of TLFre.

6

Conclusion

In this paper, we propose a novel feature reduction method for SGL via decomposition of convex sets. We also derive the set of parameter values that

lead to zero solutions of SGL. To the best of our knowledge, TLFre is the first method which is applicable to sparse models with multiple sparsity-inducing regularizers. More importantly, the proposed approach provides novel framework for developing screening methods for complex sparse models with multiple sparsity-inducing regularizers, e.g., ‘1 SVM that performs both sample and feature selection, fused Lasso and tree Lasso with more than two regularizers. Experiments on both synthetic and real data sets demonstrate the effectiveness and efficiency of TLFre. We plan to generalize the idea of TLFre to ‘1 SVM, fused Lasso and tree Lasso, which are expected to consist of multiple layers of screening.

8

2 References

- [1] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.
- [2] J. Borwein and A. Lewis. *Convex Analysis and Nonlinear Optimization*, Second Edition. Canadian Mathematical Society, 2006.
- [3] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [4] L. El Ghaoui, V. Viallon, and T. Rabhani. Safe feature elimination in sparse supervised learning. *Pacific Journal of Optimization*, 8:667?698, 2012.
- [5] J. Friedman, T. Hastie, and R. Tibshirani. A note on the group lasso and a sparse group lasso. *arXiv:1001.0736*.
- [6] N. Halko, P. Martinsson, and J. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53:217?288, 2011.
- [7] J.-B. Hiriart-Urruty. From convex optimization to nonconvex optimization. necessary and sufficient conditions for global optimality. In *Nonsmooth optimization and related topics*. Springer, 1988.
- [8] J.-B. Hiriart-Urruty. A note on the Legendre-Fenchel transform of convex composite functions. In *Nonsmooth Mechanics and Analysis*. Springer, 2006.
- [9] J. Liu, S. Ji, and J. Ye. SLEP: Sparse Learning with Efficient Projections. Arizona State University, 2009.
- [10] J. Liu and J. Ye. Moreau-Yosida regularization for grouped tree structure learning. In *Advances in neural information processing systems*, 2010.
- [11] J. Liu, Z. Zhao, J. Wang, and J. Ye. Safe screening with variational inequalities and its application to lasso. In *International Conference on Machine Learning*, 2014.
- [12] K. Ogawa, Y. Suzuki, S. Suzumura, and I. Takeuchi. Safe sample screening for Support Vector Machine. *arXiv:1401.6740*, 2014.
- [13] K. Ogawa, Y. Suzuki, and I. Takeuchi. Safe screening of non-support vectors in pathwise SVM computation. In *ICML*, 2013.
- [14] J. Peng, J. Zhu, A. Bergamaschi, W. Han, D. Noh, J. Pollack, and P. Wang. Regularized multivariate regression for indentifying master predictors with application to integrative genomics study of breast cancer. *The Annals of Applied Statistics*, 4:53?77, 2010.
- [15] A. Ruszczyński. *Nonlinear Optimization*. Princeton University Press, 2006.
- [16] N. Simon, J. Friedman., T. Hastie., and R. Tibshirani. A Sparse-Group Lasso. *Journal of Computational and Graphical Statistics*, 22:231?245, 2013.
- [17] P. Sprechmann, I. Ramirez, G. Sapiro., and Y. El-

dar. C-HiLasso: a collaborative hierarchical sparse modeling framework. *IEEE Transactions on Signal Processing*, 59:4183–4198, 2011. [18] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B*, 58:267–288, 1996. [19] R. Tibshirani, J. Bien, J. Friedman, T. Hastie, N. Simon, J. Taylor, and R. Tibshirani. Strong rules for discarding predictors in lasso-type problems. *Journal of the Royal Statistical Society Series B*, 74:245–266, 2012. [20] M. Vidyasagar. Machine learning methods in the cocomputation biology of cancer. In *Proceedings of the Royal Society A*, 2014. [21] M. Vincent and N. Hansen. Sparse group lasso and high dimensional multinomial classification. *Computational Statistics and Data Analysis*, 71:771–786, 2014. [22] J. Wang, J. Jun, and J. Ye. Efficient mixed-norm regularization: Algorithms and safe screening methods. *arXiv:1307.4156v1*. [23] J. Wang, P. Wonka, and J. Ye. Scaling svm and least absolute deviations via exact data reduction. In *International Conference on Machine Learning*, 2014. [24] J. Wang and J. Ye. Two-Layer feature reduction for sparse-group lasso via decomposition of convex sets. *arXiv:1410.4210v1*, 2014. [25] J. Wang, J. Zhou, P. Wonka, and J. Ye. Lasso screening rules via dual polytope projection. In *Advances in neural information processing systems*, 2013. [26] Z. J. Xiang and P. J. Ramadge. Fast lasso screening tests based on correlations. In *IEEE ICASSP*, 2012. [27] D. Yogatama and N. Smith. Linguistic structured sparsity in text categorization. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2014. [28] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B*, 68:49–67, 2006. [29] H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B*, 67:301–320, 2005.