
Screening Rules for Convex Problems

Anant Raj
MPI

Jakob Olbrich
ETH

Bernd Gärtner
ETH

Bernhard Schölkopf
MPI

Martin Jaggi
EPFL

Abstract

We propose a new framework for deriving screening rules for convex optimization problems. Our approach covers a large class of constrained and penalized optimization formulations, and works in two steps. First, given any approximate point, the structure of the objective function and the duality gap is used to gather information on the optimal solution. In the second step, this information is used to produce screening rules, i.e. safely identifying unimportant weight variables of the optimal solution. Our general framework leads to a large variety of useful existing as well as new screening rules for many applications. For example, we provide new screening rules for general simplex and L_1 -constrained problems, Elastic Net, squared-loss Support Vector Machines, minimum enclosing ball, as well as structured norm regularized problems, such as group lasso.

1 Introduction

¹Optimization techniques for high-dimensional problems have become the work-horses for most data-analysis and machine-learning methods. With the rapid increase of available data, major challenges occur as the number of optimization variables (weights) grows beyond capacity of current systems.

The idea of screening refers to eliminating optimization variables that are guaranteed to *not* contribute to any optimal solution, and can therefore safely be removed from problem. Such screening techniques have received increased interest in several machine learning related applications in recent years, and have been shown to lead to very significant computational efficiency improvements in various cases, in particular for many types of sparse methods. Screening techniques can be used either as a pre-processing before passing the problem to the optimizer, or also interactively during any iterative solver (called dynamic screening), to gradually reduce the problem complexity during optimization.

¹ Parts of this work have appeared in the Master's Thesis [Olbrich, 2015].

While existing screening methods were mainly relying on geometric and problem-specific properties, we in this paper take a different approach. We propose a new framework allowing screening on general convex optimization problems, using simple tools from convex duality instead of any geometric arguments. Our framework applies to a very large class of optimization problems both for constrained as well as penalized problems, including most machine learning methods of interest.

Our main contributions in this paper are summarized as follows:

1. We propose a new framework for screening for a more general class of optimization problem with a simple primal-dual structure.
2. The framework leads to a large set of new screening rules for machine learning problems that could not be screened before. Furthermore, it also recovers many existing screening rules as special cases.
3. We are able to express all screening rules using general optimization complexity notions such as smoothness or strong convexity, getting rid of problem-specific geometric properties.
4. Our proposed rules are dynamic (allowing any existing algorithm to be additionally equipped with screening) and safe (guaranteed to only eliminate truly unimportant variables).

Related Work. The concept of screening in the sense of eliminating non-influential data points to reduce the problem size has originated relatively independently in at least two communities. Coming from computational geometry, Ahipasaoğlu et al. [2008] has proposed a screening technique for the minimum enclosing ball problem for a given set of data points. Here screening can be interpreted as simply removing points which are guaranteed to lie in the strict interior of the final ball. Later Källberg and Larsson [2014] improve the threshold for this rule in the minimum enclosing ball setting.

Independently, the breakthrough work of Ghaoui et al. [2010] gave the first screening rules for the important case of sparse regression, as given in the Lasso. Since then, there

have been many extensions and alterations of the general concept. While Ghaoui et al. [2010] exploits geometric quantities to bound the the Lasso dual solution within a compact region, we recommend the survey paper by Xiang et al. [2014] for an overview of geometric methods for Lasso screening. Sphere-region based methods differ from dome-shaped regions as used in Ghaoui et al. [2010] in choosing different centers and radii to bound the dual optimal point. Apart from being geometry specific, most existing approaches such as [Wang et al., 2013, 2014, Liu et al., 2014, Ghaoui et al., 2010, Ogawa et al., 2013] are not agnostic to the regularization parameter used, but instead are restricted to perform screening along the entire regularization path (as the regularization parameter changes). This is known as sequential screening, and restricts its usability to optimization algorithms obtaining paths. In contrast, our proposed framework here allows any internal optimization algorithms to be equipped with screening.

Despite the importance of constrained problems in many applications, much less is known about screening for constrained optimization, in contrast to the case of penalized optimization problems. For the dual of the hinge loss SVM, which is a box-constrained optimization problem, Ogawa et al. [2014] proposed a geometric screening rule based on the intersection region of two spheres, in the sequential setting of varying regularizer. More recently, Zimmert et al. [2015] provided new screening rules for that case in the dynamic setting using a method similar to our approach. However their method is restricted to the SVM case.

As a first step to allow screening for more general optimization objectives, Ndiaye et al. [2015] gives duality gap based screening rules for multi-task and multi-class problems (in the penalized setting) under for a wider class of objectives f . Nevertheless, their approach is restricted to assume separability of f over the group structure, which limits the screening rules, in the sense of not covering standard group lasso for example. Also in [Shibagaki et al., 2016], authors assume the similar problem formulation as in [Ndiaye et al., 2015] but a bit more general. The focus in Shibagaki et al. [2016] is on screening rules for SVM problems rather than general framework. They derive the screening rules for SVM by considering standard hinge and ϵ -insensitive loss with regularization formulation which is close to the empirical risk minimization framework here, but has more limited applications in terms of generalization of the screening rules. We here provide screening rules for a more general framework of box constrained optimization, while hinge-loss SVM happens to be a special case of this. Our proposed approach can be shown to recover many of the other existing rules including e.g. [Ndiaye et al., 2015] and [Zimmert et al., 2015], but significantly generalizing the method to general objectives and constraints as well as regularizers.

The rest of the paper is organized as follows: In Section 2,

we discuss our framework for screening. Section 3 is devoted to deriving the information about optimal points in terms of gap functions. Sections 4 and 5 utilizes the framework and tools derived in previous sections to provide screening rules for constrained and penalized case respectively. In the end, we provide a small illustrative experiment for screening on simplex and L_1 -constrained and also discuss that which of the existing results can be recovered using our algorithm in Section 6.

2 Setup and Primal-Dual Structure

In this paper, we consider optimization problems of the following primal-dual structure. As we will see, the relationship between primal and dual objectives has many benefits, including computation of the duality gap, which allows us to have a certificate for approximation quality.

A very wide range of machine learning optimization problems can be formulated as (A) and (B), which are dual to each other:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left[\mathcal{O}_A(\mathbf{x}) := f(A\mathbf{x}) + g(\mathbf{x}) \right] \quad (\text{A})$$

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[\mathcal{O}_B(\mathbf{w}) := f^*(\mathbf{w}) + g^*(-A^\top \mathbf{w}) \right] \quad (\text{B})$$

The two problems are associated to a given data matrix $A \in \mathbb{R}^{d \times n}$, and the functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ are allowed to be arbitrary closed convex functions. The functions f^*, g^* in formulation (B) are defined as the *convex conjugates* of their corresponding counterparts f, g in (A). Here $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{w} \in \mathbb{R}^d$ are the respective variable vectors. For a given function $h : \mathbb{R}^d \rightarrow \mathbb{R}$, its conjugate is defined as

$$h^*(\mathbf{v}) := \max_{\mathbf{u} \in \mathbb{R}^d} \mathbf{v}^\top \mathbf{u} - h(\mathbf{u}).$$

The association of problems (A) and (B) is a special case of Fenchel Duality. More precisely, the relationship is called *Fenchel-Rockafellar Duality* when incorporating the linear map A as in our case, see e.g. [Borwein and Zhu, 2005, Theorem 4.4.2] or [Bauschke and Combettes, 2011, Proposition 15.18], see the Appendix A for a self-contained derivation. The two main powerful features of this general duality structure are first that it includes many more machine learning methods than more traditional duality notions, and secondly that the two problems are fully symmetric, when changing respective roles of f and g . In typical machine learning problems, the two parts typically play the roles of a data-fit (or loss) term as well as a regularization term. As we will see later, the two roles can be swapped, depending on the application.

Optimality Conditions. The first-order optimality conditions for our pair of vectors $\mathbf{w} \in \mathbb{R}^d, \mathbf{x} \in \mathbb{R}^n$ in problems (A) and (B) are given as

$$\mathbf{w} \in \partial f(A\mathbf{x}), \quad (1a) \quad -A^\top \mathbf{w} \in \partial g(\mathbf{x}), \quad (2a)$$

$$A\mathbf{x} \in \partial f^*(\mathbf{w}), \quad (1b) \quad \mathbf{x} \in \partial g^*(-A^\top \mathbf{w}) \quad (2b)$$

see e.g. [Bauschke and Combettes, 2011, Proposition 19.18]. The stated optimality conditions are equivalent to \mathbf{x}, \mathbf{w} being a saddle-point of the Lagrangian, which is given as $\mathcal{L}(\mathbf{x}, \mathbf{w}) = f^*(\mathbf{w}) - \langle A\mathbf{x}, \mathbf{w} \rangle - g(\mathbf{x})$ if $\mathbf{x} \in \text{dom}(g)$ and $\mathbf{w} \in \text{dom}(f^*)$, see Appendix A for details.

The Constrained Case. Any constrained convex optimization problem of the form

$$\min_{\mathbf{x} \in \mathcal{C}} f(A\mathbf{x}) \quad (3)$$

for a constraint set \mathcal{C} can be directly written in the form (A) by using the indicator function of the constraint set as the penalization term g . (The indicator function $\iota_{\mathcal{C}}$ of a set $\mathcal{C} \subset \mathbb{R}^n$ is defined as $\iota_{\mathcal{C}}(\mathbf{x}) := 0$ if $\mathbf{x} \in \mathcal{C}$ and $\iota_{\mathcal{C}}(\mathbf{x}) := +\infty$ otherwise.)

The Partially Separable Case. A very important special case arises when one part of the objective becomes separable. Formally, this is expressed as $g(\mathbf{x}) = \sum_{i=1}^n g_i(x_i)$ for univariate functions $g_i : \mathbb{R} \rightarrow \mathbb{R}$ for $i \in [n]$. Nicely in this case, the conjugate of g also separates as $g^*(\mathbf{y}) = \sum_i g_i^*(y_i)$. Therefore, the two optimization problems (A) and (B) write as

$$\mathcal{O}_A(\mathbf{x}) := f(A\mathbf{x}) + \sum_i g_i(x_i) \quad (\text{SA})$$

$$\mathcal{O}_B(\mathbf{w}) := f^*(\mathbf{w}) + \sum_i g_i^*(-\mathbf{a}_i^\top \mathbf{w}), \quad (\text{SB})$$

where $\mathbf{a}_i \in \mathbb{R}^d$ denotes the i -th column of A .

Crucially in this case, the optimality conditions (2a) and (2b) now become separable, that is

$$-\mathbf{a}_i^\top \mathbf{w} \in \partial g_i(x_i) \quad \forall i. \quad (4a)$$

$$x_i \in \partial g_i^*(-\mathbf{a}_i^\top \mathbf{w}) \quad \forall i. \quad (4b)$$

Note that the two other conditions (1a) and (1b) are unchanged in this case.

3 Duality Gap and Certificates

The duality gap for our problem structure provides an optimality certificate for our class of optimization problems. It will be the most important tool for us to provide guaranteed information about the optimal point (as in Section 3.2), which will then be the foundation for the second step, to perform screening on the optimal point (as we will do in the later Sections 4 and 5).

3.1 Duality Gap Structure

For the problem structure (A) and (B) as given by Fenchel-Rockafellar duality, the *duality gap* for any pair of primal

and dual variables $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{w} \in \mathbb{R}^d$ is defined as $G(\mathbf{w}, \mathbf{x}) := \mathcal{O}_A(\mathbf{x}) + \mathcal{O}_B(\mathbf{w})$. Non-negativity of the gap – that is weak duality – is satisfied by all pairs.

Most importantly, the duality gap acts as a certificate of approximation quality — the true optimum values $\mathcal{O}_A(\mathbf{x}^*)$ and $-\mathcal{O}_B(\mathbf{w}^*)$ (which are both unknown) will always lie within the (known) duality gap.

The Gap Function. For the special case of differentiable function f , we can study a simpler duality gap

$$G(\mathbf{x}) := \mathcal{O}_A(\mathbf{x}) + \mathcal{O}_B(\mathbf{w}(\mathbf{x})) \quad (5)$$

purely defined as a function of \mathbf{x} , using the optimality relation (1a), i.e. $\mathbf{w}(\mathbf{x}) := \nabla f(A\mathbf{x})$.

The Wolfe-Gap Function. For any constrained optimization problem (3) defined over a bounded set \mathcal{C} and $\mathbf{x} \in \mathcal{C}$, the Wolfe gap function (also known as Hearn gap or Frank-Wolfe gap) is defined as the difference of f to the minimum of its linearization over the same domain. Formally,

$$G_{\mathcal{C}}(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{C}} (A\mathbf{x} - A\mathbf{y})^\top \nabla f(A\mathbf{x}). \quad (6)$$

It is not hard to see that the convenient Wolfe gap function is a special case of our above defined general duality gap $G(\mathbf{x}) := \mathcal{O}_A(\mathbf{x}) + \mathcal{O}_B(\mathbf{w}(\mathbf{x}))$, for g being the indicator function of the constraint set \mathcal{C} , and $\mathbf{w}(\mathbf{x}) := \nabla f(A\mathbf{x})$. For more details, see Appendix B.1, or also [Lacoste-Julien et al., 2013, Appendix D].

3.2 Obtaining Information about the Optimal Points

As we have mentioned, any type of screening will crucially rely on first deriving safe knowledge about the unknown optimal points of our given optimization problem. Here, we will use the duality gap to obtain such knowledge on the optimal points $\mathbf{x}^* \in \mathbb{R}^n$ and $\mathbf{w}^* \in \mathbb{R}^d$ of the respective optimization problems (A) and (B) respectively. Proofs are provided in Appendix B.2.

Our first lemma shows how to bound the distance between any (feasible) current dual iterate and the solution \mathbf{w}^* using standard assumptions on the objective functions.

Lemma 1. *Consider the problem (B) with optimal solution $\mathbf{w}^* \in \mathbb{R}^d$. For f being μ -smooth, we have*

$$\|\mathbf{w} - \mathbf{w}^*\|^2 \leq \frac{2}{\mu} (f^*(\mathbf{w}) - f^*(\mathbf{w}^*)) \quad (7)$$

The following corollary will be important to derive screening rules for penalized problems in Section 5, as well as box-constrained problems (Section 4.4).

Corollary 2. *We consider the problem setup (A) and (B), and assume f is μ -smooth. Then*

$$\|\mathbf{w} - \mathbf{w}^*\|^2 \leq \frac{2}{\mu} G(\mathbf{x}). \quad (8)$$

Here $G(\mathbf{x})$ is the duality gap function as defined in equation (5).

The following two results hold for general constrained optimization problems of the form (3), where g is the indicator function of a constraint set $\mathcal{C} \subset \mathbb{R}^n$ and hence are useful for deriving screening rules for such problems.

Lemma 3. Consider problem (A) and assume that f is μ -strongly convex over a bounded set \mathcal{C} . Then it holds that

$$\|A\mathbf{x} - A\mathbf{x}^*\|_2^2 \leq \frac{1}{\mu} G_{\mathcal{C}}(\mathbf{x}), \quad (9)$$

where \mathbf{x}^* is an optimal solution and $G_{\mathcal{C}}$ is the Wolfe-Gap function of f over the bounded set \mathcal{C} .

Corollary 4. Assuming f is L -smooth as well as μ -strongly convex over a bounded set \mathcal{C} , we have

$$\|\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}^*)\| \leq \frac{L}{\sqrt{\mu}} \sqrt{G_{\mathcal{C}}(\mathbf{x})} \quad (10)$$

4 Screening Rules for Constrained Problems

In the following, we will develop screening rules for constrained optimization problems of the form (3), by exploiting the structure of the constraint set for a variety of sparsity-inducing problems. First of all, we give a general lemma which we will be using in rest of the paper to derive screening rules when any of the function in A and B is indicator function.

Lemma 5. For general constrained optimization $\min_{\mathbf{x} \in \mathcal{C}} f(A\mathbf{x})$, the optimality condition (2a) gives rise to the following optimality rule at the optimal point:

$$(A\mathbf{x}^*)^\top \mathbf{w}^* = \min_{\mathbf{z} \in \mathcal{C}} (A\mathbf{z})^\top \mathbf{w}^* \quad (11)$$

The above equation (11) also suggest that $\mathbf{x}^* = \arg \min_{\mathbf{z} \in \mathcal{C}} (A\mathbf{z})^\top \mathbf{w}^*$. Lemma 5 is very crucial in further deriving screening rules for constrained optimization problem as well as norm penalized problems whose conjugate is indicator function of the dual norm.

4.1 Simplex Constrained Problems

Optimization over unit simplex $\Delta := \{\mathbf{x} \in \mathbb{R}^n \mid x_i \geq 0, \sum_{i=1}^n x_i = 1\}$ is a important class of constrained problems (3), as it includes optimization over any finite polytope. In this case, the columns of A describe the vertices, and \mathbf{x} are barycentric coordinates representing the point $A\mathbf{x}$. Formally, $g(\mathbf{x})$ is the indicator function of the unit simplex $\mathcal{C} = \Delta$ in this case.

The following two theorems provide screening rules for simplex constrained problems. We provide all proofs in Appendix C.1.

Theorem 6. For general simplex constrained optimization $\min_{\mathbf{x} \in \Delta} f(A\mathbf{x})$, the optimality condition (2a) gives rise to the following screening rule at the optimal point, for any $i \in [n]$

$$(\mathbf{a}_i - A\mathbf{x}^*)^\top \mathbf{w}^* > 0 \Rightarrow x_i^* = 0. \quad (12)$$

In the following Theorem 7 we now assume smoothness and strong convexity of function f to provide screening rules for simplex problems, in terms of an arbitrary iterate \mathbf{x} , without knowing \mathbf{x}^* .

Theorem 7. Let f be L -smooth and μ -strongly convex over the unit simplex $\mathcal{C} = \Delta$. Then for simplex constrained optimization $\min_{\mathbf{x} \in \Delta} f(A\mathbf{x})$ we have the following screening rule, for any $i \in [n]$

$$(\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}) > L \sqrt{\frac{G_{\mathcal{C}}(\mathbf{x})}{\mu}} \|\mathbf{a}_i - A\mathbf{x}\| \Rightarrow x_i^* = 0. \quad (13)$$

Our general screening rules for simplex constrained problems as in Theorem 7 allows many practical implications. For example, new screening rules for squared loss SVM and minimum enclosing ball problem come as a direct consequence.

Squared Hinge Loss SVM. The squared hinge-loss SVM problem in its dual form is formulated as

$$\min_{\mathbf{x} \in \Delta} [f(A\mathbf{x}) := \frac{1}{2} \mathbf{x}^\top A^\top A \mathbf{x}] \quad (14)$$

over a unit simplex constraint $\mathbf{x} \in \Delta \subset \mathbb{R}^n$. Here for given data examples $\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_n \in \mathbb{R}^d$ and corresponding labels $y_i \in \pm 1$, the matrix A collects the columns $\mathbf{a}_i = y_i \bar{\mathbf{a}}_i$, see e.g. Tsang et al. [2005]. We obtain the following novel screening rule for square loss SVM:

Corollary 8. For the squared hinge loss SVM (14) we have the screening rule

$$\begin{aligned} (\mathbf{a}_i - A\mathbf{x})^\top A\mathbf{x} &> \sqrt{\max_i (A\mathbf{x} - \mathbf{a}_i)^\top A\mathbf{x}} \|\mathbf{a}_i - A\mathbf{x}\| \\ &\Rightarrow x_i^* = 0. \end{aligned} \quad (15)$$

Minimum Enclosing Ball. The primal and dual for the minimum enclosing ball problem is given as the following pair of optimization formulations (16) and (17) respectively.

$$\min_{\mathbf{c} \in \mathbb{R}^d, r \in \mathbb{R}} r^2 \quad \text{s.t.} \quad \|\mathbf{c} - \mathbf{a}_i\|_2^2 \leq r^2 \quad \forall i \in [n] \quad (16)$$

$$\min_{\mathbf{x} \in \Delta \subset \mathbb{R}^n} \mathbf{x}^\top A^\top A \mathbf{x} + \mathbf{c}^\top \mathbf{x}, \quad (17)$$

where \mathbf{c} is a vector whose i^{th} element c_i is $-\mathbf{a}_i^\top \mathbf{a}_i$, see for example [Matoušek and Gärtner, 2007] or our Appendix C.1.

Corollary 9. For the minimum enclosing ball problem (16) we have the screening rule

$$\begin{aligned} & (\mathbf{e}_i - \mathbf{x})^\top (2A^\top A\mathbf{x} + \mathbf{c}) > \\ & 2\sqrt{\frac{1}{2} \max_i (\mathbf{x} - \mathbf{e}_i)^\top (2A^\top A\mathbf{x} + \mathbf{c})} \|\mathbf{a}_i - A\mathbf{x}\| \\ \Rightarrow x_i^* &= 0. \end{aligned} \quad (18)$$

Our result improves upon the known rules by Källberg and Larsson [2014], Ahipasaoğlu et al. [2008] by providing a broader selection criterion (18).

4.2 L_1 -Constrained Problems

L_1 -constrained formulations are very widely used in order to induce sparsity in the variables. Here below we provide results for screening on general L_1 -constrained problems, that is $\min_{\mathbf{x} \in \mathcal{C}} f(A\mathbf{x})$ for $\mathcal{C} = L_1 \subset \mathbb{R}^n$ (or a scaled version of the L_1 -ball). Proofs are provided in Appendix C.2.

Theorem 10. For general L_1 -constrained optimization $\min_{\mathbf{x} \in L_1} f(A\mathbf{x})$, the optimality condition (2a) gives rise to the following screening rule at the optimal point, for any $i \in [n]$

$$|\mathbf{a}_i^\top \mathbf{w}^*| + (A\mathbf{x}^*)^\top \mathbf{w}^* < 0 \Rightarrow x_i^* = 0. \quad (19)$$

Using only a current iterate \mathbf{x} instead of an optimal point, we obtain screening for general smooth and strongly convex function f :

Theorem 11. Let f be L -smooth and μ -strongly convex over the L_1 -ball. Then for L_1 -constrained optimization $\min_{\mathbf{x} \in L_1} f(A\mathbf{x})$ we have the following screening rule, for any $i \in [n]$

$$\begin{aligned} & |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) \\ & + L(\|\mathbf{a}_i\|_2 + \|A\mathbf{x}\|_2) \sqrt{\frac{G_{\mathcal{C}}(\mathbf{x})}{\mu}} < 0 \\ \Rightarrow x_i^* &= 0 \end{aligned} \quad (20)$$

4.3 Elastic Net Constrained Problems

Elastic net regularization as an alternative to L_1 is often used in practice, and can outperform the Lasso, while still enjoying a similar sparsity of representation Zou and Hastie [2005]. The elastic net is given by the expression

$$\alpha \|\mathbf{x}\|_1 + \frac{(1-\alpha)}{2} \|\mathbf{x}\|_2^2.$$

Here below we provide novel result for screening on general elastic net constrained problems, that is $\min_{\mathbf{x} \in \mathcal{C}} f(A\mathbf{x})$ for \mathcal{C} being the elastic net constraint, or a scaled version of it. Proofs are provided in Appendix C.3.

Theorem 12. For general elastic net constrained optimization $\min_{\mathbf{x} \in L_E} f(A\mathbf{x})$ where $L_E := \{\mathbf{x} \in \mathbb{R}^n \mid \alpha \|\mathbf{x}\|_1 + \frac{(1-\alpha)}{2} \|\mathbf{x}\|_2^2 \leq 1\}$, the optimality condition (2a) gives rise to the following screening rule at the optimal point, for any $i \in [n]$

$$|\mathbf{a}_i^\top \mathbf{w}^*| + (A\mathbf{x}^*)^\top \mathbf{w}^* \left[\frac{\alpha}{1 + \frac{(1-\alpha)}{2} \|\mathbf{x}^*\|_2^2} \right] < 0 \Rightarrow x_i^* = 0$$

Using only a current iterate \mathbf{x} instead of an optimal point, we obtain screening for general smooth and strongly convex function f :

Theorem 13. Let f be L -smooth and μ -strongly convex over the elastic net norm ball. Then for elastic net constrained optimization $\min_{\mathbf{x} \in L_E} f(A\mathbf{x})$ we have the following screening rule, for any $i \in [n]$

$$\begin{aligned} & |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) \left[\frac{2\alpha}{3-\alpha} \right] \\ & + L(\|\mathbf{a}_i\|_2 + \|A\mathbf{x}\|_2 \left[\frac{2\alpha}{3-\alpha} \right]) \sqrt{\frac{G_{\mathcal{C}}(\mathbf{x})}{\mu}} < 0 \\ \Rightarrow x_i^* &= 0 \end{aligned} \quad (21)$$

Note that both above results also recover the L_1 constrained case as a special case, when $\alpha \rightarrow 1$.

4.4 Screening for Box Constrained Problems

Box-constrained problems are important in several machine learning applications, including SVMs. After variable rescaling, w.l.o.g. we can assume the constraint set $\mathcal{C} = \square := \{\mathbf{x} \in \mathbb{R}^n \mid 0 \leq x_i \leq 1\}$. We derive screening rules for predicting both if a variable will take the upper or lower constraint.

Theorem 14. Let f be L -smooth. Then for box-constrained optimization $\min_{\mathbf{x} \in \square} f(A\mathbf{x})$, we obtain the following screening rules, for any $i \in [n]$

$$\begin{aligned} & \mathbf{a}_i^\top \nabla f(A\mathbf{x}) - \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} > 0 \Rightarrow x_i^* = 0, \text{ and} \\ & \mathbf{a}_i^\top \nabla f(A\mathbf{x}) + \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} < 0 \Rightarrow x_i^* = 1. \end{aligned}$$

Box constrained optimization problems arise very often in machine learning problem. Hinge loss SVM happens to one of many special cases of box-constrained optimization problem.

Hinge Loss SVM. The dual of the classical support vector machine with hinge loss, when not using a bias value, is a box-constrained problem. As a direct consequence of Theorem 14 we therefore obtain screening rules for SVM with hinge loss and no bias. The primal formulation of the SVM in this setting, for a regularization parameter $C > 0$,

is

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \boldsymbol{\varepsilon} \in \mathbb{R}^n} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \mathbf{1}^\top \boldsymbol{\varepsilon} \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{a}_i \geq 1 - \varepsilon_i \quad \forall i \in [n] \\ & \varepsilon_i \geq 0 \quad \forall i \in [n] \end{aligned} \quad (22)$$

Corollary 15. *For SVM with hinge loss and no bias as given in (22), we have the screening rules*

$$\begin{aligned} \mathbf{a}_i^\top A \mathbf{x} - \|\mathbf{a}_i\|_2 \sqrt{2G(\mathbf{x})} > 0 &\Rightarrow x_i^* = 0, \text{ and} \\ \mathbf{a}_i^\top A \mathbf{x} + \|\mathbf{a}_i\|_2 \sqrt{2G(\mathbf{x})} < 0 &\Rightarrow x_i^* = C. \end{aligned}$$

where $\mathbf{x} \in \mathbb{R}^n$ is any feasible dual point.

We get similar screening rules for hinge loss SVM as in [Zimmert et al., 2015] as well as in [Shibagaki et al., 2016]. The closest known result to our Corollary 15 for screening in hinge loss SVM is given in Zimmert et al. [2015] and [Shibagaki et al., 2016]. The work of Zimmert et al. [2015] also covers the kernelized SVM case, and improves the threshold given in our Corollary 15 by a constant of $\sqrt{2}$. In Appendix C.4, we show that our more general approach here can also be adjusted to gain this constant factor.

5 Screening for Penalized Problems

In this section we will develop screening methods for general penalized convex optimization problems of the form (A) and (B). The cornerstone application are L_1 regularized problems, for which we now develop screening rules with general cost function f . We show in Appendix D.1 that our method can reproduce the screening rules of Ndiaye et al. [2015] as special cases, whereas their method does not directly extend to general f . Beyond L_1 problems, we also describe new screening rules for elastic net regularized problems, as well as the important case of structured norm regularized optimization.

5.1 L_1 -Penalized Problems

The next theorem describes a screening rule for general L_1 -penalized problems, under a smoothness assumption on function f . Proofs for are given in Appendix D.1.

Theorem 16. *Consider an L_1 -regularized optimization problem of the form*

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(A\mathbf{x}) + \lambda \|\mathbf{x}\|_1 \quad (23)$$

If f is L -smooth, then the following screening rule holds for all $i \in [n]$:

$$|\mathbf{a}_i^\top \nabla f(A\mathbf{x})| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{2L G(\mathbf{x})} \Rightarrow x_i^* = 0$$

By careful observation of the expression in Theorem 16, it is easy to find a connection between our screening rule and the geometric sphere test method based screening Xiang

et al. [2014]. The general idea behind the sphere test is to consider the maximum value of the objective function in a spherical region which contains the optimal dual variable. We discuss this connection in more detail in Appendix D.3.

Also, in Appendix D.1, we discuss the special cases of squared loss regression and logistic loss regression with L_1 penalization. These results are presented in Corollaries 24 and 25 as direct consequences of Theorem 16. Both of the corollaries can also be derived from the framework discussed in the paper Ndiaye et al. [2015].

5.2 Elastic-Net Penalized Problems

In the next corollary, we present a novel screening rule for the elastic net squared loss regression problem.

Corollary 17. *Consider the elastic net regression formulation*

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda_2 \|\mathbf{x}\|_2^2 + \lambda_1 \|\mathbf{x}\|_1 \quad (24)$$

The following screening rule holds for all $i \in [n]$:

$$\begin{aligned} |(\mathbf{a}_i^\top A + 2\lambda_2 \mathbf{e}_i^\top) \mathbf{x} - \mathbf{a}_i^\top \mathbf{b}| < \\ \lambda_1 - \sqrt{2(\mathbf{a}_i^\top \mathbf{a}_i + 2\lambda_2)G(\mathbf{x})} \Rightarrow x_i^* = 0. \end{aligned}$$

We also recover existing screening rules for elastic net regularized problem with more general objective f using our frameworks, in Appendix D.1, see Lemma 26 and Theorem 27 which has been earlier derived in [Shibagaki et al., 2016]. In the proof, we derive screening rules from both the formulation (SA) and (SB) using optimality condition (4a) and (4b) which is novel as well as help us to understand the property useful in deriving screening rules for elastic net penalized problems.

5.3 Structured Norm Penalized Problems

Here in this section we present screening rules for non-overlapping group norm regularized problems. Group-norm regularization is widely used to induce sparsity in terms of groups of variables of the the solution of the optimization problem. The most prominent example is the group lasso (ℓ_2/ℓ_1 -regularization). Here in this section we mostly discuss screening for general objectives with an ℓ_2/ℓ_1 -regularization. Proofs are provided in Appendix D.2.

Group Norm - ℓ_2/ℓ_1 Regularization. In the following, we use the notation $\{\mathbf{x}_1 \cdots \mathbf{x}_G\}$ to express a vector \mathbf{x} as a partition of the groups of variables, such that $\mathbf{x}^\top = [\mathbf{x}_1^\top, \mathbf{x}_2^\top \cdots \mathbf{x}_G^\top]$. Correspondingly, the matrix A can be denoted as the concatenation of the respective columns $A = [A_1 \ A_2 \cdots A_G]$.

Theorem 18. *For ℓ_2/ℓ_1 -regularized optimization problem of the form*

$$\min_{\mathbf{x}} f(A\mathbf{x}) + \sum_{g=1}^G \sqrt{\rho_g} \|\mathbf{x}_g\|_2$$

Assuming f is L -smooth, then the following screening rule holds for all groups g :

$$\|A_g^\top \nabla f(A\mathbf{x})\|_2 + \sqrt{2L} \|A_g\|_{Fro} < \sqrt{\rho_g} \Rightarrow \mathbf{x}_g^* = 0.$$

Corollary 19. Group Lasso Regression with Squared Loss - For the group lasso formulation

$$\min_{\mathbf{x}} \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \sum_{g=1}^G \sqrt{\rho_g} \|\mathbf{x}_g\|_2$$

we have the following screening rule for all groups g :

$$\|A_g^\top (A\mathbf{x} - \mathbf{b})\|_2 + \sqrt{2G(\mathbf{x})} \|A_g\|_{Fro} < \lambda \sqrt{\rho_g} \Rightarrow \mathbf{x}_g^* = 0.$$

Group lasso regression is widely used in applications as an working example case of structured norm penalization. The framework of Ndiaye et al. [2015] does not directly provide screening rules for the group lasso, due to the fact that they require f to be partially separable over the groups as well as special structural requirement of the formulation, in contrast to our more general Theorem 18. Similarly, Lee and Xing [2014] is also restricted to least-squares f objective.

6 Illustrative Experiments

While the contribution of our paper is on the theoretical generality and the collection of new screening applications, we will still briefly illustrate the performance of some of the proposed screening algorithms, for the classical examples of simplex constrained and L_1 -constrained problems. We compare the fraction of active variables and the Wolfe-Gap function as optimization algorithm progress.

We consider the optimization problem of the form $\min_{\mathbf{x} \in \mathcal{B}_{L_1}} \|A\mathbf{x} - \mathbf{b}\|_2^2$. \mathcal{B}_{L_1} is a scaled L_1 -ball with radius

35. $A \in \mathbb{R}^{3000 \times 600}$ is a random Gaussian matrix and a noisy measurement $\mathbf{b} = A\mathbf{x}^*$ where \mathbf{x}^* is a sparse vector of +1 and -1 with only 70 non zeros entries. We solve the above optimization problem using the Frank-Wolfe algorithm (pair-wise variant, see Lacoste-Julien and Jaggi [2015]). Before putting this optimization problem into the solver we convert this problem into the barycentric representation which is $\min_{\mathbf{x}_\Delta \in \Delta} \|A_\Delta \mathbf{x}_\Delta - \mathbf{b}\|_2^2$. The relation between the transformed variable and original variable can be given by $A_\Delta = [A \mid -A]$ and $\mathbf{x} = [I_n \mid -I_n] \mathbf{x}_\Delta$. For more details see [Jaggi, 2014].

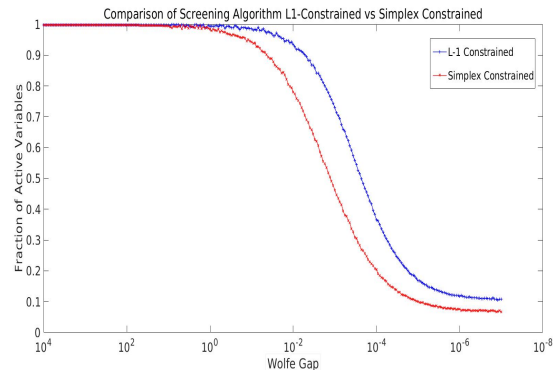


Figure 1: Simplex- vs L_1 -constrained Screening

Dataset/ No. of Samples	No Screening (Simplex)	Screening (Simplex)
Synth1 5000	13.1 sec	11.7sec
Synth2 10000	28.3 sec	23.1 sec
RCV1 20242	18.6 min	13.5 min
news20B 19996	33.4 min	25.2 min

Table 1: Simplex-constrained screening, clock time

Dataset/ No. of Samples	No Screening (ℓ_1 -constr.)	Screening (ℓ_1 -constr.)
Synth1 5000	13.1	12.2 sec
Synth2 10000	28.3 sec	24.7 sec
RCV1 20242	18.6 min	14.9 min
news20B 19996	33.4 min	27.1 min

Table 2: L_1 -constrained screening, clock time

Now we apply our Theorems 11 and 7 on variable of \mathbf{x} and \mathbf{x}_Δ respectively to screen, in order to compare the two alternative screening approaches on the same problem. Note that the Wolfe gap is identical in both parameterizations, for any \mathbf{x} . One important point to note here is that dimension of \mathbf{x}_Δ is the double of the dimension of \mathbf{x} , and any L_1 -coordinate value x_i is zero if and only if both “duplicate” variables $x_{\Delta,i}$ and $x_{\Delta,n+i}$ are zero, where n is the dimensionality of \mathbf{x} .

Therefore, the simplex variant (with more variables) performs a more fine-grained variant of screening, where we can screen each of the sign patterns separately for each variable. In Fig 1, the blue curve illustrates the screening efficiency for the L_1 -constrained screening case, while the red curve illustrate simplex constrained screening. Our theorems 11 and 7 are well in line with the phenomena in Fig 1. For the L_1 -constrained case, the screening starts relatively at later stage than simplex case due to the fact that in Equation (21), two out of three terms are absolute values of some quantity and hence it is very tough to compensate both of them by the third quantity, in order for the entire sum to become negative. Hence in the beginning this rule can often be ineffective. As algorithm progresses, the duality gap becomes smaller and screening starts but at the same time

the gradient (and therefore gap) also starts to decay which brings the trade-off shown in the plot. For both variants, screening becomes slow towards the end.

We also report the time taken to reach a duality gap of 10^{-7} with both the approaches mentioned above (simplex constrained and L_1 -constrained) on for different datasets. The first two datasets (*Synth1* and *Synth2*) are generated under the same setting described earlier but *Synth1* with 5000 samples and *Synth2* with 10000 samples. *RCVI* is a real world dataset having 20,242 samples and 47,236 data dimensions. *news20Binary* is also a real world dataset having 19,996 entries and 1,355,191 dimensions. Below in Tables 1 and 2, we describe the running time of the optimization methods to reach a duality gap threshold of 10^{-7} with or without screening. On *RCVI* dataset we try the feature learning with L_1 -norm ball constraint of 200 and on *news20Binary* we use L_1 -norm ball constraint of 35. In the case of *RCVI* and *news20Binary*, A is the data matrix and \mathbf{b} is the label of each instance in the dataset. From Tables 1 and 2 it is also evident that simplex screening rule is more tighter than the L_1 -constrained screening rule.

7 Discussion

We have presented a unified way to derive screening rules for general constrained and penalized optimization problems. For both cases, our framework crucially utilizes the structure of piece-wise linearity of the problem at hand. For the constrained case, we showed that screening rules follow from the piece-wise linearity of the boundary of the constraint set.

The crucial property is that at non-differentiable boundary points, the normal cone – i.e. the sub-differential of the indicator function of the constraint set – becomes a relatively large set. Under moderate assumptions on the objective function, we are able to guarantee that also the gradient of an optimal point must lie in this same cone region, leading to screening.

On the other hand for penalized optimization problems, we are able to derive screening rules from either piece-wise linearity of the penalty function, or as well from exploiting piece-wise linearity of the constraint set arising from the dual (conjugate) of the penalty function.

Acknowledgements. The authors would like to thank Julia Wysling, Rohit Babbar and Peter Bühlmann for valuable discussions.

References

- S. D. Ahipasaoğlu, P. Sun, and M. Todd. Linear Convergence of a Modified Frank–Wolfe Algorithm for Computing Minimum-Volume Enclosing Ellipsoids. *Optimization Methods and Software*, 23(1):5–19, Feb. 2008.
- H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. CMS Books in Mathematics. Springer New York, New York, NY, 2011.
- J. M. Borwein and Q. Zhu. *Techniques of Variational Analysis and Nonlinear Optimization*. Canadian Mathematical Society Books in Math, Springer New York, 2005.
- S. P. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- L. E. Ghaoui, V. Viallon, and T. Rabbani. Safe Feature Elimination for the LASSO and Sparse Supervised Learning Problems. *arXiv.org*, Sept. 2010.
- M. Jaggi. An Equivalence between the Lasso and Support Vector Machines. In *Regularization, Optimization, Kernels, and Support Vector Machines*, pages 1–26. Chapman and Hall/CRC, Oct. 2014.
- S. M. Kakade, S. Shalev-Shwartz, and A. Tewari. On the duality of strong convexity and strong smoothness: Learning applications and matrix regularization. Technical report, Toyota Technological Institute - Chicago, USA, 2009.
- L. Källberg and T. Larsson. Improved Pruning of Large Data Sets for the Minimum Enclosing Ball Problem. *Graphical Models*, July 2014.
- S. Lacoste-Julien and M. Jaggi. On the Global Linear Convergence of Frank-Wolfe Optimization Variants. In *NIPS 2015 - Advances in Neural Information Processing Systems 28*, 2015.
- S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-Coordinate Frank-Wolfe Optimization for Structural SVMs. In *ICML 2013 - Proceedings of the 30th International Conference on Machine Learning*, 2013.
- S. Lee and E. P. Xing. Screening Rules for Overlapping Group Lasso. *arXiv*, Oct. 2014.
- J. Liu, Z. Zhao, J. Wang, and J. Ye. Safe Screening with Variational Inequalities and Its Application to Lasso. In *ICML 2014 - Proceedings of the 31st International Conference on Machine Learning*, pages 289–297, 2014.
- J. Matoušek and B. Gärtner. *Understanding and using linear programming*. Springer, 2007.
- E. Ndiaye, O. Fercoq, A. Gramfort, and J. Salmon. GAP Safe screening rules for sparse multi-task and multi-class models. In *NIPS 2015 - Advances in Neural Information Processing Systems 28*, pages 811–819, 2015.
- K. Ogawa, Y. Suzuki, and I. Takeuchi. Safe Screening of Non-Support Vectors in Pathwise SVM Computation. In *ICML*, pages 1382–1390, 2013.
- K. Ogawa, Y. Suzuki, S. Suzumura, and I. Takeuchi. Safe Sample Screening for Support Vector Machines. *arXiv.org*, Jan. 2014.
- J. Olbrich. Screening Rules for Convex Problems. Master’s thesis, ETH Zürich, 2015.
- A. Shibagaki, M. Karasuyama, K. Hatano, and I. Takeuchi. Simultaneous Safe Screening of Features and Samples in Doubly Sparse Modeling. In *ICML 2016 - Proceedings of the 33th International Conference on Machine Learning*, pages 1577–1586, 2016.
- V. Smith, S. Forte, M. Jaggi, and M. I. Jordan. L_1 -Regularized Distributed Optimization: A Communication-Efficient Primal-Dual Framework. *arXiv cs.LG*, 2015.
- I. W. Tsang, J. T. Kwok, and P.-M. Cheung. Core Vector Machines: Fast SVM Training on Very Large Data Sets. *Journal of Machine Learning Research*, 6:363–392, 2005.
- J. Wang, B. Lin, P. Gong, P. Wonka, and J. Ye. Lasso Screening Rules via Dual Polytope Projection. In *NIPS 2014 - Advances in Neural Information Processing Systems 27*, 2013.
- J. Wang, J. Zhou, J. Liu, P. Wonka, and J. Ye. A Safe Screening Rule for Sparse Logistic Regression. In *NIPS 2014 - Advances in Neural Information Processing Systems 27*, pages 1053–1061, 2014.
- Z. J. Xiang, Y. Wang, and P. J. Ramadge. Screening Tests for Lasso Problems. *arXiv.org*, May 2014.
- J. Zimmert, C. S. de Witt, G. Kerg, and M. Kloft. Safe screening for support vector machines. In *NIPS Workshop on Optimization for Machine Learning*, pages 1–5, Dec. 2015.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.

A Primal Dual Structure (Section 2)

The relation of our primal and dual problems (A) and (B) is standard in convex analysis, and is a special case of the concept of Fenchel Duality. Using the combination with the linear map A as in our case, the relationship is called *Fenchel-Rockafellar Duality*, see e.g. [Borwein and Zhu, 2005, Theorem 4.4.2] or [Bauschke and Combettes, 2011, Proposition 15.18]. For completeness, we here illustrate this correspondence with a self-contained derivation of the duality.

Proof. Starting with the original formulation (A), we introduce a helper variable vector $\mathbf{v} \in \mathbb{R}^d$ representing $\mathbf{v} = A\boldsymbol{\alpha}$. Then optimization problem (A) becomes:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} f(\mathbf{v}) + g(\boldsymbol{\alpha}) \quad \text{such that } \mathbf{v} = A\boldsymbol{\alpha}. \quad (25)$$

Introducing Lagrange multipliers $\mathbf{w} \in \mathbb{R}^d$, the Lagrangian is given by:

$$L(\boldsymbol{\alpha}, \mathbf{v}; \mathbf{w}) := f(\mathbf{v}) + g(\boldsymbol{\alpha}) + \mathbf{w}^\top (A\boldsymbol{\alpha} - \mathbf{v}).$$

The dual problem of (A) follows by taking the infimum with respect to both $\boldsymbol{\alpha}$ and \mathbf{v} :

$$\begin{aligned} \inf_{\boldsymbol{\alpha}, \mathbf{v}} L(\boldsymbol{\alpha}, \mathbf{v}) &= \inf_{\mathbf{v}} \{f(\mathbf{v}) - \mathbf{w}^\top \mathbf{v}\} + \inf_{\boldsymbol{\alpha}} \{g(\boldsymbol{\alpha}) + \mathbf{w}^\top A\boldsymbol{\alpha}\} \\ &= -\sup_{\mathbf{v}} \{\mathbf{w}^\top \mathbf{v} - f(\mathbf{v})\} - \sup_{\boldsymbol{\alpha}} \{(-\mathbf{w}^\top A)\boldsymbol{\alpha} - g(\boldsymbol{\alpha})\} \end{aligned} \quad (26)$$

$$= -f^*(\mathbf{w}) - g^*(-A^\top \mathbf{w}). \quad (27)$$

We change signs and turn the maximization of the dual problem (27) into a minimization and thus we arrive at the dual formulation (B) as claimed:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[\mathcal{O}_B(\mathbf{w}) := f^*(\mathbf{w}) + g^*(-A^\top \mathbf{w}) \right].$$

The Partially Separable Case. For $g(\mathbf{x})$ is separable, i.e. $g(\mathbf{x}) = \sum_{i=1}^n g_i(x_i)$ for univariate functions $g_i : \mathbb{R} \rightarrow \mathbb{R}$ for $i \in [n]$, the primal-dual structure remains the separable. In this case, the conjugate of g also separates as $g^*(\mathbf{y}) = \sum_i g_i^*(y_i)$. Therefore, in terms of the the primal-dual structure (A) and (B) we obtain the separable special case (SA) and (SB). \square

Optimality Conditions. The first-order optimality conditions follow from the standard definition of the conjugate functions in the Fenchel dual problem, see also e.g. Borwein and Zhu [2005], Bauschke and Combettes [2011].

Proof. The first-order optimality conditions for our pair of vectors $\mathbf{w} \in \mathbb{R}^d$, $\mathbf{x} \in \mathbb{R}^n$ in problems (A) and (B) are given by equations (1a), (2a), (1b) and (2b). The proof directly comes from equation (26) by separately writing optimizing conditions for two expressions $\mathbf{w}^\top \mathbf{v} - f(\mathbf{v})$ and $(-\mathbf{w}^\top A)\boldsymbol{\alpha} - g(\boldsymbol{\alpha})$ in equation (26).

Crucially in the partially separable case, the optimality conditions (2a) and (2b) become separable. Comparing the expressions (SA) and (A), we see that $g(\mathbf{x}) = \sum_i g_i(x_i)$ and hence

$$g^*(\mathbf{x}) = \sum_i g_i^*(x_i)$$

Hence by applying (2a) and (2b) we obtain the separable optimality conditions (4a) and (4b). \square

B Duality Gap and Objective Function Properties

B.1 Wolfe Gap as a Special Case of Duality Gap

Proof. To see this as a special case of general duality gap of the problem formulation, we consider the constraint as indicator function of set \mathcal{C} such that $g(\mathbf{x}) = \iota_{\mathcal{C}}(\mathbf{x})$. Now from the definition of the Wolfe gap function

$$G_{\mathcal{C}}(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{C}} (A\mathbf{x} - A\mathbf{y})^\top \partial f(A\mathbf{x})$$

Here $\partial f(A\mathbf{x})$ is an arbitrary subgradient of f at the candidate position \mathbf{x} , and $\iota_{\mathcal{C}}^*(\mathbf{y}) := \sup_{\mathbf{s} \in \mathcal{C}} \langle \mathbf{s}, \mathbf{y} \rangle$ is the support function of \mathcal{C} . Now writing the general duality gap $G(\mathbf{x})$ as

$$\begin{aligned} G(\mathbf{x}) &:= \mathcal{O}_A(\mathbf{x}) + \mathcal{O}_B(\mathbf{w}(\mathbf{x})) \\ &:= f(A\mathbf{x}) + \iota_{\mathcal{C}}(\mathbf{x}) + f^*(\mathbf{w}(\mathbf{x})) + \iota_{\mathcal{C}}^*(-(A^\top \mathbf{w}(\mathbf{x}))) \end{aligned}$$

the last term disappears since we assumed $\mathbf{x} \in \mathcal{C}$. Using the definition of the Fenchel conjugate, one has the Fenchel-Young inequality, i.e.

$$f^*(\mathbf{w}) := \max_{\mathbf{u} \in \mathbb{R}^d} \mathbf{w}^\top \mathbf{u} - f(\mathbf{u}) \Rightarrow f^*(\mathbf{w}) + f(\mathbf{u}) \geq \mathbf{w}^\top \mathbf{u}$$

The above holds with equality if \mathbf{w} is chosen as a subgradient of f at $\mathbf{u} = A\mathbf{x}$. Therefore, using our first-order optimality mapping $\mathbf{w}(\mathbf{x}) := \partial f(A\mathbf{x})$, we have

$$G(\mathbf{x}) = (A\mathbf{x})^\top \partial f(A\mathbf{x}) + \iota_{\mathcal{C}}^*(-(A^\top \mathbf{w}(\mathbf{x}))) = G_{\mathcal{C}}(\mathbf{x})$$

This derivation is adapted from [Lacoste-Julien et al., 2013, Appendix D]. □

B.2 Obtaining Information about the Optimal Points

Lemma 20 (Conjugates of Indicator Functions and Norms).

- i) The conjugate of the indicator function $\iota_{\mathcal{C}}$ of a set $\mathcal{C} \subset \mathbb{R}^n$ (not necessarily convex) is the support function of the set \mathcal{C} , that is $\iota_{\mathcal{C}}^*(\mathbf{x}) = \sup_{\mathbf{s} \in \mathcal{C}} \langle \mathbf{s}, \mathbf{x} \rangle$
- ii) The conjugate of a norm is the indicator function of the unit ball of the dual norm.

Proof. [Boyd and Vandenberghe, 2004, Example 3.24 and 3.26] □

Lemma 21. Assume that f is a closed and convex function then f^* is μ -strongly convex with respect to a norm $\|\cdot\|$ if and only if f is $1/\mu$ -Lipschitz gradient with respect to dual norm $\|\cdot\|_*$.

Proof. [Kakade et al., 2009, Theorem 3] □

Proof of Lemma 1. From the definition of μ -strongly convex function, we know that

$$\begin{aligned} f^*(\mathbf{w}) &\geq f^*(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^\top \nabla f^*(\mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\geq f^*(\mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \end{aligned}$$

The first inequality follows directly by using the first order optimality condition for \mathbf{w}^* being optimal. For any optimal point \mathbf{w}^* and another feasible point \mathbf{w} ,

$$(\mathbf{w} - \mathbf{w}^*)^\top \nabla f^*(\mathbf{w}^*) \geq 0.$$

Hence, $\|\mathbf{w}^* - \mathbf{w}\|_2^2 \leq \frac{2}{\mu} (f^*(\mathbf{w}) - f^*(\mathbf{w}^*))$ □

Proof of Corollary 2. This statement directly comes from (1) and the definition of the duality gap. By definition we know that the true optimum values $\mathcal{O}_A(\mathbf{x}^*)$ and $-\mathcal{O}_B(\mathbf{w}^*)$ respectively for primal (A) and dual formulation (B) will always lie within the duality gap which implies

$$G(\mathbf{x}) \geq \mathcal{O}_B(\mathbf{w}) - \mathcal{O}_B(\mathbf{w}^*)$$

By equation (B), we know that $\mathcal{O}_B(\mathbf{w}) = f^*(\mathbf{w}) + g^*(-A^\top \mathbf{w}^*)$

Now since f^* is μ -strongly convex function and g^* is convex hence,

$$f^*(\mathbf{w}) \geq f^*(\mathbf{w}^*) + \nabla f^*(\mathbf{w}^*)^\top (\mathbf{w} - \mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \quad (28)$$

$$g^*(-A^\top \mathbf{w}) \geq g^*(-A^\top \mathbf{w}^*) + \nabla g^*(-A^\top \mathbf{w}^*)^\top (-A^\top \mathbf{w} + A^\top \mathbf{w}^*) \quad (29)$$

Hence by adding equation (28) and (29), we get

$$\begin{aligned}\mathcal{O}_B(\mathbf{w}) &\geq \mathcal{O}_B(\mathbf{w}^*) + (\nabla f^*(\mathbf{w}^*) - A\nabla g^*(-A^\top \mathbf{w}^*))^\top (\mathbf{w} - \mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\Rightarrow \mathcal{O}_B(\mathbf{w}) \geq \mathcal{O}_B(\mathbf{w}^*) + \nabla \mathcal{O}_B(\mathbf{w}^*)^\top (\mathbf{w} - \mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2\end{aligned}$$

At optimal point \mathbf{w}^* , $\nabla \mathcal{O}_B(\mathbf{w}^*)^\top (\mathbf{w} - \mathbf{w}^*) \geq 0$.

Hence,

$$G(\mathbf{x}) \geq \mathcal{O}_B(\mathbf{w}) - \mathcal{O}_B(\mathbf{w}^*) \geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

□

Proof of Lemma 3. From the definition of μ -strong convexity of f and using optimality condition,

$$\mu \|A\mathbf{x} - A\mathbf{x}^*\|^2 \leq (A\mathbf{x} - A\mathbf{x}^*)^\top (\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}^*)) \quad (30)$$

$$\leq (A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}) \quad (31)$$

$$\leq G_{\mathcal{C}}(\mathbf{x}) \quad (32)$$

Equation (30) comes from the definition of μ -strong convexity.

Equation (31) is first order optimality condition for \mathbf{x}^* being optimal which implies

$$(A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \geq 0$$

The inequality (32) follows by the definition of the gap function given in (6). □

Proof of Corollary 4. This comes by definition of L -smooth functions and Lemma 3. From the definition,

$$\begin{aligned}\|\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}^*)\| &\leq L \|A\mathbf{x} - A\mathbf{x}^*\| \\ &\leq \frac{L}{\sqrt{\mu}} \sqrt{G_{\mathcal{C}}(\mathbf{x})}\end{aligned}$$

Second inequality directly comes from Lemma 3. □

C Screening on Constrained Problems

Lemma 22. Let \mathcal{C} be a convex set, and $\iota_{\mathcal{C}}$ be its indicator function, then

1. For $\mathbf{x} \notin \mathcal{C}$, $\partial \iota_{\mathcal{C}}(\mathbf{x}) = \emptyset$
2. For $\mathbf{x} \in \mathcal{C}$, we have that $\mathbf{w} \in \partial \iota_{\mathcal{C}}(\mathbf{x})$ if $\mathbf{w}^\top (\mathbf{z} - \mathbf{x}) \leq 0 \quad \forall \mathbf{z} \in \mathcal{C}$

Proof. Let $\mathcal{C} \subseteq \mathbb{R}^n$ be a closed convex set. Then subgradient of indicator function $\iota_{\mathcal{C}}(\mathbf{x})$ at \mathbf{x} will be vectors \mathbf{u} which satisfy

$$\begin{aligned}\iota_{\mathcal{C}}(\mathbf{z}) &\geq \iota_{\mathcal{C}}(\mathbf{x}) + \mathbf{u}^\top (\mathbf{z} - \mathbf{x}) \quad \forall \mathbf{z} \in \text{dom}(\iota_{\mathcal{C}}) \\ &\Rightarrow \iota_{\mathcal{C}}(\mathbf{z}) \geq \iota_{\mathcal{C}}(\mathbf{x}) + \mathbf{u}^\top (\mathbf{z} - \mathbf{x}) \quad \forall \mathbf{z} \in \mathbb{R}^n\end{aligned} \quad (33)$$

If $\text{int}(\mathcal{C})$ represents the interior of the set \mathcal{C} such that it contains n -dimensional ball of radius $r > 0$, and $Bd(\mathcal{C})$ represents boundary of the set \mathcal{C} . Now we have to assume various cases for proving Lemma 22.

Case 1 We evaluate Equation (33) when $\mathbf{x} \in \text{int}(\mathcal{C})$. Equation (33) becomes

$$\iota_{\mathcal{C}}(\mathbf{z}) \geq \mathbf{u}^\top (\mathbf{z} - \mathbf{x}) \quad \forall \mathbf{z} \in \mathbb{R}^n$$

Now since the above equation is satisfied for all $\mathbf{z} \in \mathbb{R}^n$, we assume $\mathbf{z} \in \text{int}(\mathcal{C})$ such that $(\mathbf{z} - \mathbf{x})$ can be anywhere in the ball. Hence \mathbf{u} needs to be 0 in this case.

Case 2 In this case we assume $x \in Bd(\mathcal{C})$. That gives

$$\iota_{\mathcal{C}}(z) \geq \mathbf{u}^\top(z - \mathbf{x}) \quad \forall z \in \mathbb{R}^n$$

If we take $z \in \mathcal{C}$ then \mathbf{u} satisfies $\mathbf{u}^\top(z - \mathbf{x}) \leq 0 \quad \forall z \in \mathcal{C}$

If $z \notin \mathcal{C}$ then \mathbf{u} can take all the value. Hence taking intersection, \mathbf{u} satisfies

$$\mathbf{u}^\top(z - \mathbf{x}) \leq 0 \quad \forall z \in \mathcal{C}$$

Case 3 When we assume $x \notin \mathcal{C}$, we get

$$\iota_{\mathcal{C}}(z) \geq +\infty + \mathbf{u}^\top(z - \mathbf{x}) \quad \forall z \in \mathbb{R}^n$$

If we again take $z \in \mathcal{C}$ then no finite \mathbf{u} can satisfy the equation $\iota_{\mathcal{C}}(z) \geq +\infty + \mathbf{u}^\top(z - \mathbf{x}) \quad \forall z \in \mathcal{C}$ because $\iota_{\mathcal{C}}(z) = 0$ if $z \in \mathcal{C}$.

And if $z \notin \mathcal{C} \Rightarrow \iota_{\mathcal{C}}(z) = +\infty$ then again nothing can be said about the vector \mathbf{u} . Hence by convention it is assumed that $x \notin \mathcal{C} \Rightarrow \mathbf{u} \in \emptyset$

By the above arguments we conclude that,

1. For $x \notin \mathcal{C}$, $\partial \iota_{\mathcal{C}}(x) = \emptyset$
2. For $x \in \mathcal{C}$, we have that $\mathbf{w} \in \partial \iota_{\mathcal{C}}(x)$ if $\mathbf{w}^\top(z - \mathbf{x}) \leq 0 \quad \forall z \in \mathcal{C}$

Hence the claim made in Lemma 22 is proved. □

Proof of Lemma 5. From Lemma 22, we know the expression for subgradient of the indication function $\iota_{\mathcal{C}}$

$$\begin{aligned} \partial g(\mathbf{x}^*) &= \{ \mathbf{s} \mid \forall z \in \mathcal{C} \quad \mathbf{s}^\top(z - \mathbf{x}^*) \leq 0 \} \\ &= \{ \mathbf{s} \mid \forall z \in \mathcal{C} \quad \mathbf{s}^\top z \leq \mathbf{s}^\top \mathbf{x}^* \} \end{aligned} \quad (34)$$

Now, by the optimality condition (2a), $-A^\top \mathbf{w}^* \in \partial g(\mathbf{x}^*)$ and since this holds, hence $-A^\top \mathbf{w}^*$ should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g(\mathbf{x}^*)$ according to conditions in equation (39). Hence,

$$(-A^\top \mathbf{w}^*)^\top z \leq (-A^\top \mathbf{w}^*)^\top \mathbf{x}^* \quad \forall z \in \mathcal{C} \quad (35)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq (A^\top \mathbf{w}^*)^\top z \quad \forall z \in \mathcal{C} \quad (36)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq \min_z (A^\top \mathbf{w}^*)^\top z \quad s.t. \quad z \in \mathcal{C} \quad (37)$$

$$\Rightarrow (A\mathbf{x}^*)^\top \mathbf{w}^* \leq \min_{z \in \mathcal{C}} (Az)^\top \mathbf{w}^* \quad s.t. \quad z \in \mathcal{C} \quad (38)$$

Since \mathbf{x}^* is a feasible point hence $(A\mathbf{x}^*)^\top \mathbf{w}^* = \min_{z \in \mathcal{C}} (Az)^\top \mathbf{w}^* \quad s.t. \quad \mathbf{x}^*, z \in \mathcal{C}$. □

C.1 Screening on Simplex Constrained Problems (Section 4.1)

General Simplex Constrained Screening

Proof of Theorem 6. In the simplex case, we have $g(\mathbf{x}) = \iota_{\Delta}(\mathbf{x})$ and by Lemma 22

$$\begin{aligned} \partial g(\mathbf{x}^*) &= \{ \mathbf{s} \mid \forall z \in \Delta \quad \mathbf{s}^\top(z - \mathbf{x}^*) \leq 0 \} \\ &= \{ \mathbf{s} \mid \forall z \in \Delta \quad \mathbf{s}^\top z \leq \mathbf{s}^\top \mathbf{x}^* \} \end{aligned} \quad (39)$$

Now, by the optimality condition (2a), $-A^\top \mathbf{w}^* \in \partial g(\mathbf{x}^*)$ and since this holds, hence $-A^\top \mathbf{w}^*$ should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g(\mathbf{x}^*)$ according to conditions in equation (39). Hence,

$$(-A^\top \mathbf{w}^*)^\top \mathbf{z} \leq (-A^\top \mathbf{w}^*)^\top \mathbf{x}^* \quad \forall \mathbf{z} \in \Delta \quad (40)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad \forall \mathbf{z} \in \Delta \quad (41)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq \min_z (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad s.t. \quad \mathbf{z} \in \Delta \quad (42)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq \min_i \mathbf{a}_i^\top \mathbf{w}^* \quad (43)$$

$$\Rightarrow (A\mathbf{x}^*)^\top \mathbf{w}^* \leq \min_i \mathbf{a}_i^\top \mathbf{w}^* \quad (44)$$

$$\Rightarrow (A\mathbf{x}^*)^\top \mathbf{w}^* = \min_i \mathbf{a}_i^\top \mathbf{w}^* \quad (45)$$

Equation (43) is due to the fact that \mathbf{z} lie in the simplex, hence minimum value of $(A^\top \mathbf{w}^*)^\top \mathbf{z}$ is $\min_i \mathbf{a}_i^\top \mathbf{w}^*$ and equation (45) also comes from the same fact that \mathbf{x}^* lie in the simplex and hence $(A\mathbf{x}^*)^\top \mathbf{w}^*$ can not be smaller than $\min_i \mathbf{a}_i^\top \mathbf{w}^*$. That implies these two quantities need to be equal and all the i 's where this equality doesn't hold refers to $x_i^* = 0$ for all such i 's.

$$\mathbf{a}_i^\top \mathbf{w}^* > (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \Rightarrow x_i = 0$$

$$(\mathbf{a}_i - A\mathbf{x}^*)^\top \mathbf{w}^* > 0 \Rightarrow x_i = 0$$

□

Proof of Theorem 7. From the optimality condition (1a), we have $\mathbf{w}^* = \nabla f(A\mathbf{x}^*)$ since f is differentiable. Hence,

$$(\mathbf{a}_i - A\mathbf{x}^*)^\top \mathbf{w}^* = (\mathbf{a}_i - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \quad (46)$$

$$= (\mathbf{a}_i - A\mathbf{x}^* + A\mathbf{x} - A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \quad (47)$$

$$= (\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) + (A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \quad (48)$$

$$\geq (\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \quad \{\text{From the optimality of } f(A\mathbf{x})\} \quad (49)$$

$$= (\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}) - (\mathbf{a}_i - A\mathbf{x})^\top (\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}^*)) \quad (50)$$

$$\geq (\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}) - \|\mathbf{a}_i - A\mathbf{x}\| \|\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}^*)\| \quad (51)$$

$$\geq (\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}) - L \sqrt{\frac{G_C(\mathbf{x})}{\mu}} \|\mathbf{a}_i - A\mathbf{x}\| \quad (52)$$

Eq. (49) comes from the fact that at the optimal point \mathbf{x}^* , the inequality $(A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \geq 0$ holds $\forall \mathbf{x}$. Equation (52) comes from Corollary 4 for smooth function f over a constrained set \mathcal{C} .

Hence from Theorem 6, we obtain the screening rule

$$(\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x}) > L \sqrt{\frac{G_C(\mathbf{x})}{\mu}} \|\mathbf{a}_i - A\mathbf{x}\| \Rightarrow x_i^* = 0$$

□

Screening for Squared Hinge Loss SVM.

Proof of Corollary 8. Theorem 7 is directly applicable to problems of the form (14). The objective function $f(\mathbf{y}) = f(A\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top A^\top A \mathbf{x}$ is strongly convex with parameter $\mu = 1$. Also the derivative ∇f is Lipschitz-continuous with parameter $L = 1$. To obtain an upper bound on the distance between any approximate solution and the optimal solution $\|A\mathbf{x} - A\mathbf{x}^*\|$, we employ Lemma 3. Since the constrained of the optimization problem is unit simplex and hence the value of Wolfe gap function $G_C(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{C}} (A\mathbf{x} - A\mathbf{y})^\top \nabla f(A\mathbf{x})$ as defined in Section 3 will be attained on one of the vertices. So, $G_C(\mathbf{x}) = \max_{i \in 1 \dots m} (A\mathbf{x} - \mathbf{a}_i)^\top A\mathbf{x}$. Finally, Theorem 7 gives us the screening rule for squared hinge loss SVM:

$$(\mathbf{a}_i - A\mathbf{x})^\top A\mathbf{x} > \sqrt{\max_{i \in 1 \dots m} (A\mathbf{x} - \mathbf{a}_i)^\top A\mathbf{x}} \|\mathbf{a}_i - A\mathbf{x}\| \Rightarrow x_i^* = 0 \quad (53)$$

□

Screening on Minimum Enclosing Ball.

Minimum Enclosing Ball - Given a set of n points, \mathbf{a}_1 to \mathbf{a}_n in \mathbb{R}^d , the minimum enclosing ball is defined as the smallest ball $B_{c,r}$ with center c and radius r , i.e.: $B_{c,r} := \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{c} - \mathbf{x}\| \leq r\}$, such that all points \mathbf{a}_i lie in its interior. In this set-up, screening means to identify points \mathbf{a}_i lying in the interior of the optimal ball B_{c^*,r^*} . Removing those points from the problem does not change the optimal ball.

Proof of Corollary 9. The minimum enclosing ball problem can be formulated as an optimization problem of the form given in Equation (16):

$$\min_{c,r} r^2 \quad \text{s.t. } \|\mathbf{c} - \mathbf{a}_i\|_2^2 \leq r^2 \quad \forall i \in [n]$$

As we have seen, the dual formulation can be written in the form of Equation (17) as given in [Matoušek and Gärtner, 2007, Chapter 8.7]:

$$\min_{\mathbf{x}} \mathbf{x}^\top A^\top A \mathbf{x} - \sum_{j=1}^p \mathbf{a}_j^\top \mathbf{a}_j x_j \quad \text{s.t. } \mathbf{x} \in \Delta$$

Now the function $\mathbf{x}^\top A^\top A \mathbf{x} - \sum_{j=1}^p \mathbf{a}_j^\top \mathbf{a}_j x_j$ is strongly convex in $A\mathbf{x}$ with parameter $\mu = 2$. Since the constrained of the optimization problem is unit simplex and hence the value of the Wolfe gap function $G_C(\mathbf{x}) := \max_{\mathbf{y} \in C} (A\mathbf{x} - A\mathbf{y})^\top \nabla f(A\mathbf{x})$ as defined in Section 3 will be attained at one of the vertices of unit simplex. Hence Corollary 4 gives $G_C(\mathbf{x}) = \sqrt{\frac{1}{2} \max_i (\mathbf{x} - \mathbf{e}_i)^\top (2A^\top A \mathbf{x} + \mathbf{c}')}$. Now applying the findings of Theorem 7, we get a sufficient condition for \mathbf{a}_i to be non-influential, i.e. \mathbf{a}_i lies in the interior of the MEB. But before that we will simplify the left hand side of the theorem 7 a bit. $(\mathbf{a}_i - A\mathbf{x})^\top \nabla f(A\mathbf{x})$ can be written as $(\mathbf{e}_i - \mathbf{x})^\top A^\top \nabla f(A\mathbf{x})$. Hence we get our result claimed in Corollary 9.

$$(\mathbf{e}_i - \mathbf{x})^\top (2A^\top A \mathbf{x} + \mathbf{c}') > 2 \sqrt{\frac{1}{2} \max_j (\mathbf{x} - \mathbf{e}_j)^\top (2A^\top A \mathbf{x} + \mathbf{c}')} \|\mathbf{a}_i - A\mathbf{x}\| \Rightarrow x_i^* = 0 \quad (54)$$

That means \mathbf{a}_i is non influential. □

C.2 Screening on L_1 -ball Constrained Problems

Proof of Theorem 10. In the constrained Lasso case, we have $g(\mathbf{x}) = \nu_{\mathcal{B}_{L_1}}(\mathbf{x})$ and by Lemma 22

$$\begin{aligned} \partial g(\mathbf{x}^*) &= \{ \mathbf{s} \mid \forall \mathbf{z} \in \mathcal{B}_{L_1} \quad \mathbf{s}^\top (\mathbf{z} - \mathbf{x}^*) \leq 0 \} \\ &= \{ \mathbf{s} \mid \forall \mathbf{z} \in \mathcal{B}_{L_1} \quad \mathbf{s}^\top \mathbf{z} \leq \mathbf{s}^\top \mathbf{x}^* \} \end{aligned} \quad (55)$$

Now, by the optimality condition (2a), $-A^\top \mathbf{w}^* \in \partial g(\mathbf{x}^*)$ and since this holds, hence $-A^\top \mathbf{w}^*$ should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g(\mathbf{x}^*)$ according to conditions in equation (70). Hence,

$$(-A^\top \mathbf{w}^*)^\top \mathbf{z} \leq (-A^\top \mathbf{w}^*)^\top \mathbf{x}^* \quad \forall \mathbf{z} \in \mathcal{B}_{L_1} \quad (56)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad \forall \mathbf{z} \in \mathcal{B}_{L_1} \quad (57)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq \min_z (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad \text{s.t. } \mathbf{z} \in \mathcal{B}_{L_1} \quad (58)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq -\max_i |\mathbf{a}_i^\top \mathbf{w}^*| \quad (59)$$

$$\Rightarrow (A\mathbf{x}^*)^\top \mathbf{w}^* \leq -\max_i |\mathbf{a}_i^\top \mathbf{w}^*| \quad (60)$$

$$\Rightarrow (A\mathbf{x}^*)^\top \mathbf{w}^* = -\max_i |\mathbf{a}_i^\top \mathbf{w}^*| \quad (61)$$

Equation (59) is due to the fact that \mathbf{z} lie in the L_1 -ball and hence minimum value of $(A^\top \mathbf{w}^*)^\top \mathbf{z}$ is $-\max_i |\mathbf{a}_i^\top \mathbf{w}^*|$ and Equation (61) also comes from the same fact that \mathbf{x}^* lie in the L_1 -ball and hence $(A\mathbf{x}^*)^\top \mathbf{w}^*$ can not be smaller than $-\max_i |\mathbf{a}_i^\top \mathbf{w}^*|$. That implies these two quantities need to be equal and all the i 's where this equality doesn't hold refers to $x_i^* = 0$ for all such i 's. Hence whenever these two quantities are not equal this holds:

$$\begin{aligned} -|\mathbf{a}_i^\top \mathbf{w}^*| &> (A\mathbf{x}^*)^\top \mathbf{w}^* \Rightarrow x_i^* = 0 \\ \Rightarrow |\mathbf{a}_i^\top \mathbf{w}^*| + (A\mathbf{x}^*)^\top \mathbf{w}^* &< 0 \Rightarrow x_i^* = 0 \end{aligned}$$

□

Proof of Theorem 11. Using optimality condition (1a), we know that $\mathbf{w}^* \in \partial f(A\mathbf{x})$

$$|\mathbf{a}_i^\top \mathbf{w}^*| + (A\mathbf{x}^*)^\top \mathbf{w}^* = |\mathbf{a}_i^\top \nabla f(A\mathbf{x}^*)| + (A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \quad (62)$$

$$= |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}) + \nabla f(A\mathbf{x}^*))| + (A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \quad (63)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x}^* - A\mathbf{x} + A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \quad (64)$$

$$= |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) - (A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \quad (65)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \quad (66)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}) + \nabla f(A\mathbf{x})) \quad (67)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) + (A\mathbf{x})^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x})) \quad (68)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) + L(\|\mathbf{a}_i\| + \|A\mathbf{x}\|) \sqrt{\frac{G_C(\mathbf{x})}{\mu}} \quad (69)$$

Eq. (65) comes from the fact that at the optimal point \mathbf{x}^* , the inequality $(A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \geq 0$ holds $\forall \mathbf{x}$. Hence using Theorem 10, Lemma 3 and Corollary 4, we get the screening rule for L_1 constrained as whenever,

$$|\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) + L(\|\mathbf{a}_i\| + \|A\mathbf{x}\|) \sqrt{\frac{G_C(\mathbf{x})}{\mu}} < 0 \Rightarrow \mathbf{x}_i^* = 0 \quad \square$$

C.3 Screening on Elastic Net Constrained Problems

Proof of Theorem 12. Formulation :

$$\begin{aligned} & \min_{\mathbf{x}} f(A\mathbf{x}) \\ & \text{s.t. } \alpha \|\mathbf{x}\|_1 + \frac{(1-\alpha)}{2} \|\mathbf{x}\|_2^2 \leq 1 \\ & \Rightarrow \alpha \sum_{i=1}^n |\mathbf{x}_i| + \frac{(1-\alpha)}{2} \sum_{i=1}^n \mathbf{x}_i^2 \leq 1 \end{aligned}$$

In the elastic net constrained case, we have $g(\mathbf{x}) = \iota_{\mathcal{B}_{L_E}}(\mathbf{x})$ where $\iota_{\mathcal{B}_{L_E}}$ is elastic net norm ball. That implies

$$\mathbf{x} \in \mathcal{B}_{L_E} : \alpha \|\mathbf{x}\|_1 + (1-\alpha) \|\mathbf{x}\|_2^2 \leq 1$$

From the subgradient of indicator function and optimality condition for A and B framework

$$\begin{aligned} \partial g(\mathbf{x}^*) &= \{ \mathbf{s} \mid \forall \mathbf{z} \in \mathcal{B}_{L_1} \quad \mathbf{s}^\top (\mathbf{z} - \mathbf{x}^*) \leq 0 \} \\ &= \{ \mathbf{s} \mid \forall \mathbf{z} \in \mathcal{B}_{L_1} \quad \mathbf{s}^\top \mathbf{z} \leq \mathbf{s}^\top \mathbf{x}^* \} \end{aligned} \quad (70)$$

Now, by the optimality condition (2a), $-A^\top \mathbf{w}^* \in \partial g(\mathbf{x}^*)$ and since this holds, hence $-A^\top \mathbf{w}^*$ should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g(\mathbf{x}^*)$ according to conditions in equation (70). Hence,

$$(-A^\top \mathbf{w}^*)^\top \mathbf{z} \leq (-A^\top \mathbf{w}^*)^\top \mathbf{x}^* \quad \forall \mathbf{z} \in \mathcal{B}_{L_E} \quad (71)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad \forall \mathbf{z} \in \mathcal{B}_{L_E} \quad (72)$$

$$\Rightarrow (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \leq \min_{\mathbf{z}} (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad \text{s.t. } \mathbf{z} \in \mathcal{B}_{L_E} \quad (73)$$

Since \mathbf{x}^* is a feasible point hence $(A^\top \mathbf{w}^*)^\top \mathbf{x}^* = \min_{\mathbf{z}} (A^\top \mathbf{w}^*)^\top \mathbf{z} \quad \text{s.t. } \mathbf{x}^*, \mathbf{z} \in \mathcal{B}_{L_E}$. At the point where above equality

hold \mathbf{x}^* would be same as optimal \mathbf{z} . Hence the problem reduces to,

$$\begin{aligned} & \min (A^\top \mathbf{w}^*)^\top \mathbf{z} \\ & s.t \ \alpha \|\mathbf{z}\|_1 + \frac{(1-\alpha)}{2} \|\mathbf{z}\|_2^2 \leq 1 \\ & \Rightarrow \alpha \sum_{i=1}^n |z_i| + \frac{(1-\alpha)}{2} \sum_{i=1}^n z_i^2 \leq 1 \end{aligned}$$

Without the loss of generality let us assume that for $i \in \{1 \dots m\}$, $z_i \geq 0$ and $i \in \{m+1 \dots n\}$, $z_i \leq 0$. Hence the optimization problem can be written as :

$$\begin{aligned} & \min (A^\top \mathbf{w}^*)^\top \mathbf{z} \\ & s.t \ \alpha \left(\sum_{i=1}^m z_i - \sum_{i=m+1}^n z_i \right) + \frac{(1-\alpha)}{2} \sum_{i=1}^n z_i^2 \leq 1 \\ & \quad -z_i \leq 0 \text{ for } i \in \{1 \dots m\} \\ & \quad z_i \leq 0 \text{ for } i \in \{m+1 \dots n\} \end{aligned} \tag{74}$$

Writing lagrangian for optimization problem (74)

$$\mathcal{L}(\mathbf{z}, \lambda, u) = (A^\top \mathbf{w}^*)^\top \mathbf{z} - \sum_{i=1}^m \lambda_i z_i + \sum_{i=m+1}^n \lambda_i z_i + u \left(\alpha \left(\sum_{i=1}^m z_i - \sum_{i=m+1}^n z_i \right) + \frac{(1-\alpha)}{2} \sum_{i=1}^n z_i^2 - 1 \right)$$

Also optimization conditions are $\lambda_i \geq 0$, $\lambda_i z_i = 0$ and $\alpha \left(\sum_{i=1}^m z_i - \sum_{i=m+1}^n z_i \right) + (1-\alpha) \sum_{i=1}^n z_i^2 = 1$. Also we conclude from above that if $\lambda_i > 0 \Rightarrow z_i = 0$. From first order optimality condition, For $i \in \{1 \dots m\}$

$$\mathbf{a}_i^\top \mathbf{w}^* - \lambda_i = -u(\alpha + (1-\alpha)|z_i|) \tag{75}$$

For $i \in \{m+1 \dots n\}$

$$\mathbf{a}_i^\top \mathbf{w}^* + \lambda_i = -u(\alpha + (1-\alpha)|z_i|) \tag{76}$$

Now in equations (75) and (76) we multiply by z_i and add them. We get:

$$(A^\top \mathbf{w}^*)^\top \mathbf{z} + u \left[1 + \frac{(1-\alpha)}{2} \|\mathbf{z}\|_2^2 \right] = 0 \tag{77}$$

From equations (75), (76), (77) and optimality conditions discussed above we get:

$$|\mathbf{a}_i^\top \mathbf{w}^*| + (A^\top \mathbf{w}^*)^\top \mathbf{z} \left[\frac{\alpha + (1-\alpha)|z_i|}{1 + \frac{(1-\alpha)}{2} \|\mathbf{z}\|_2^2} \right] < 0 \Rightarrow z_i = 0$$

As discussed above \mathbf{x}^* share same solution as optimal \mathbf{z} . Hence

$$|\mathbf{a}_i^\top \mathbf{w}^*| + (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \left[\frac{\alpha}{1 + \frac{(1-\alpha)}{2} \|\mathbf{x}^*\|_2^2} \right] < 0 \Rightarrow x_i^* = 0$$

□

Proof of Theorem 13. Using optimality condition (1a), we know that $\mathbf{w}^* \in \partial f(A\mathbf{x})$

$$|\mathbf{a}_i^\top \mathbf{w}^*| + (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \left[\frac{\alpha}{1 + \frac{(1-\alpha)}{2} \|\mathbf{x}^*\|_2^2} \right] \leq |\mathbf{a}_i^\top \mathbf{w}^*| + (A^\top \mathbf{w}^*)^\top \mathbf{x}^* \left[\frac{2\alpha}{3-\alpha} \right] \quad (78)$$

$$= |\mathbf{a}_i^\top \nabla f(A\mathbf{x}^*)| + (A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \left[\frac{2\alpha}{3-\alpha} \right] \quad (79)$$

$$= |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}) - \nabla f(A\mathbf{x}) + \nabla f(A\mathbf{x}^*))| + (A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \left[\frac{2\alpha}{3-\alpha} \right] \quad (80)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x}^* - A\mathbf{x} + A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \left[\frac{2\alpha}{3-\alpha} \right] \quad (81)$$

$$= |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \left[\frac{2\alpha}{3-\alpha} \right] - (A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \left[\frac{2\alpha}{3-\alpha} \right] \quad (82)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}^*) \left[\frac{2\alpha}{3-\alpha} \right] \quad (83)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}) + \nabla f(A\mathbf{x})) \left[\frac{2\alpha}{3-\alpha} \right] \quad (84)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + |\mathbf{a}_i^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x}))| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) \left[\frac{2\alpha}{3-\alpha} \right] + (A\mathbf{x})^\top (\nabla f(A\mathbf{x}^*) - \nabla f(A\mathbf{x})) \left[\frac{2\alpha}{3-\alpha} \right] \quad (85)$$

$$\leq |\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) \left[\frac{2\alpha}{3-\alpha} \right] + L(\|\mathbf{a}_i\| + \|A\mathbf{x}\| \left[\frac{2\alpha}{3-\alpha} \right]) \sqrt{\frac{G_C(\mathbf{x})}{\mu}} \quad (86)$$

Eq. (82) comes from the fact that at the optimal point \mathbf{x}^* , the inequality $(A\mathbf{x} - A\mathbf{x}^*)^\top \nabla f(A\mathbf{x}^*) \geq 0$ holds $\forall \mathbf{x}$. Hence using Theorem 10, Lemma 3 and Corollary 4, we get the screening rule for L_1 constrained as whenever,

$$|\mathbf{a}_i^\top \nabla f(A\mathbf{x})| + (A\mathbf{x})^\top \nabla f(A\mathbf{x}) \left[\frac{2\alpha}{3-\alpha} \right] + L(\|\mathbf{a}_i\|_2 + \|A\mathbf{x}\|_2 \left[\frac{2\alpha}{3-\alpha} \right]) \sqrt{\frac{G_C(\mathbf{x})}{\mu}} < 0 \Rightarrow \mathbf{x}_i^* = 0$$

□

C.4 Screening for Box Constrained Problems

Screening for General Box Constrained Problems (Section 4.4)

Proof of Theorem 14. The box-constrained case can be seen in the form of the partially separable optimization problem pair (SA) and (SB). According to optimality condition (4a) for this case, we have

$$-\mathbf{a}_i^\top \mathbf{w}^* \in \partial g_i(x_i^*) \quad \forall i \quad (87)$$

Now from the definition of subgradient for an indicator function as given in Lemma 22. Also since x_i is a number now, we will get rid of the transpose here.

$$\begin{aligned} \partial g(x_i^*) &= \{s \mid 0 \leq z \leq C, \quad s(z - x_i^*) \leq 0 \} \\ &= \{s \mid 0 \leq z \leq C, \quad sz \leq sx_i^* \} \end{aligned} \quad (88)$$

Now, by the optimality condition (4a), $-\mathbf{a}_i^\top \mathbf{w}^* \in \partial g(x_i^*)$ and since this holds, hence $-\mathbf{a}_i^\top \mathbf{w}^*$ should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g(x_i^*)$ according to conditions in Equation (88). Hence,

$$\begin{aligned} (-\mathbf{a}_i^\top \mathbf{w}^*)z &\leq (-\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \quad \forall z \text{ s.t. } 0 \leq z \leq C, \\ \Rightarrow \min_z (\mathbf{a}_i^\top \mathbf{w}^*)z &\geq (\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \quad \text{s.t. } 0 \leq z \leq C \end{aligned} \quad (89)$$

Now (89) can be manipulated in two ways

Case 1

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{w}^* > 0 &\Rightarrow \min_z (\mathbf{a}_i^\top \mathbf{w}^*)z \geq (\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \quad \text{s.t. } 0 \leq z \leq C \\ &\Rightarrow 0 \geq (\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \end{aligned}$$

But since $\mathbf{a}_i^\top \mathbf{w}^* > 0$ and also $x_i^* \geq 0$ hence $(\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \not\leq 0$. This implies $(\mathbf{a}_i^\top \mathbf{w}^*)x_i^* = 0$ and hence if $\mathbf{a}_i^\top \mathbf{w}^* > 0 \Rightarrow x_i^* = 0$

Case 2

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{w}^* < 0 &\Rightarrow \min_z (\mathbf{a}_i^\top \mathbf{w}^*)z \geq (\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \quad \text{s.t. } 0 \leq z \leq C \\ &\Rightarrow (\mathbf{a}_i^\top \mathbf{w}^*)C \geq (\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \end{aligned}$$

But since $\mathbf{a}_i^\top \mathbf{w}^* < 0$ and also $x_i^* \leq C$ hence $(\mathbf{a}_i^\top \mathbf{w}^*)x_i^* \not\leq (\mathbf{a}_i^\top \mathbf{w}^*)C$. This implies $(\mathbf{a}_i^\top \mathbf{w}^*)x_i^* = (\mathbf{a}_i^\top \mathbf{w}^*)C$ and hence if $\mathbf{a}_i^\top \mathbf{w}^* < 0 \Rightarrow x_i^* = C$

Final optimality arguments can be given as

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{w}^* > 0 &\Rightarrow x_i^* = 0 \\ \mathbf{a}_i^\top \mathbf{w}^* < 0 &\Rightarrow x_i^* = C \end{aligned} \quad (90)$$

Now

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{w}^* &= \mathbf{a}_i^\top (\mathbf{w}^* + \mathbf{w} - \mathbf{w}) = \mathbf{a}_i^\top \mathbf{w} + \mathbf{a}_i^\top (\mathbf{w}^* - \mathbf{w}) \\ \mathbf{a}_i^\top \mathbf{w} - \|\mathbf{a}_i\|_2 \|\mathbf{w} - \mathbf{w}^*\|_2 &\leq \mathbf{a}_i^\top \mathbf{w}^* \leq \mathbf{a}_i^\top \mathbf{w} + \|\mathbf{a}_i\|_2 \|\mathbf{w} - \mathbf{w}^*\|_2 \end{aligned} \quad (91)$$

Since f is L -Lipschitz gradient hence f^* is $1/L$ -strongly convex, hence using Lemmas 1 and 21, Equation (90) becomes

$$\mathbf{a}_i^\top \mathbf{w} - \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \leq \mathbf{a}_i^\top \mathbf{w}^* \leq \mathbf{a}_i^\top \mathbf{w} + \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \quad (92)$$

Hence using equation (92) and earlier arguments we get,

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{w}^* > 0 &\Rightarrow x_i^* = 0 \\ \Rightarrow \mathbf{a}_i^\top \mathbf{w} - \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} &> 0 \Rightarrow x_i^* = 0 \end{aligned}$$

And if

$$\begin{aligned} \mathbf{a}_i^\top \mathbf{w}^* < 0 &\Rightarrow x_i^* = C \\ \Rightarrow \mathbf{a}_i^\top \mathbf{w} + \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} &> 0 \Rightarrow x_i^* = C \end{aligned}$$

□

Screening on SVM with hinge loss and no bias

Proof of Corollary 15. Here the primal problem is given by:

$$\begin{aligned} \min_{\mathbf{w}, \epsilon} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \mathbf{1}^\top \epsilon \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{a}_i \geq 1 - \epsilon_i \quad \forall i \in \{1 : p\} \\ & \epsilon_i \geq 0 \quad \forall i \in \{1 : p\} \end{aligned} \quad (93)$$

A dual formulation of the problem can be written as:

$$\begin{aligned} \min_{\mathbf{x}} \quad & -\mathbf{x}^\top \mathbf{1} + \frac{1}{2} \mathbf{x}^\top A^\top A \mathbf{x} \\ \text{s.t.} \quad & 0 \leq \mathbf{x} \leq C \mathbf{1} \end{aligned} \quad (94)$$

Theorem 14 is applied on the dual formulation. The objective function $\frac{1}{2} \mathbf{x}^\top A^\top A \mathbf{x} - \mathbf{x}^\top \mathbf{1}$ is strongly convex with parameter 1 and its derivative Lipschitz continuous with parameter 1. The duality gap between primal and dual feasible points $G(\mathbf{w}, \epsilon, \mathbf{x})$ is now used as suboptimality certificate which can play the role of the upper bound $\|\mathbf{w} - \mathbf{w}^*\|$ using Lemma 2. For a given \mathbf{x} a primal feasible point can be obtained by setting $\mathbf{w} = A\mathbf{x}$ and ϵ minimal such that the first constraint of the primal problem is satisfied. Using the obtained point for the duality gap, it only depends on the point \mathbf{x} . All together this gives the screening rule:

$$\mathbf{a}_i^\top A \mathbf{x} + 1 > \|\mathbf{a}_i\| \sqrt{2G(\mathbf{x})} \Rightarrow x_i^* = 0 \quad (95)$$

$$\mathbf{a}_i^\top A \mathbf{x} + 1 < -\|\mathbf{a}_i\| \sqrt{2G(\mathbf{x})} \Rightarrow x_i^* = C \quad (96)$$

□

Note - Since the primal and dual of hinge loss SVM have very nice structure with smooth quadratic function with an addition to piece-wise linear convex function, hence it is not hard to show that both primal and dual function is 1 strongly convex as shown in Zimmert et al. [2015]. For more detailed proof, we recommend to go through Zimmert et al. [2015]. Now for an instance, if we write duality gap function as a function of \mathbf{w} then

$$G(\mathbf{w}) \geq G(\mathbf{w}^*) + \nabla G(\mathbf{w}^*)^\top (\mathbf{w} - \mathbf{w}^*) + \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

Since strong duality hold in SVM case, hence at optimal point \mathbf{w}^* , $G(\mathbf{w}^*) = 0$. Finally we get,

$$G(\mathbf{w}) \geq \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

Hence the screening rule comes out as given in Zimmert et al. [2015]:

$$\mathbf{a}_i^\top A \mathbf{x} + 1 > \|\mathbf{a}_i\| \sqrt{G(\mathbf{x})} \Rightarrow x_i^* = 0 \quad (97)$$

$$\mathbf{a}_i^\top A \mathbf{x} + 1 < -\|\mathbf{a}_i\| \sqrt{G(\mathbf{x})} \Rightarrow x_i^* = C \quad (98)$$

D Screening on Penalized Problems

D.1 Screening L_1 -regularized Problems

Lemma 23. *Considering general L_1 -regularized optimization problems*

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(A\mathbf{x}) + \lambda \|\mathbf{x}\|_1 \quad (99)$$

At optimum points \mathbf{x}^* and dual optimal point \mathbf{w}^* , the following rule is satisfied for the above problem formulation (99) :

$$|\mathbf{a}_i^\top \mathbf{w}^*| < \lambda \Rightarrow x_i^* = 0$$

Proof. Since the optimization problem (99) comes under the partially separable framework and we can use the first order optimality condition (4a) as well as (4b) to derive screening rules for the problem. Also we know that, the conjugate of the norm function is the indicator function of its dual norm ball. By the optimality condition (4b), we know that

$$x_i \in \partial g_i^*(-\mathbf{a}_i^\top \mathbf{w})$$

here g_i^* is the indicator function written as $\iota_{L_\infty}(-\mathbf{a}_i^\top \mathbf{w})$. Hence for the indicator function g^* by Lemma 22

$$\begin{aligned} \partial g_i^*(-\mathbf{a}_i^\top \mathbf{w}^*) &= \left\{ s \mid \forall \mathbf{z} \text{ s.t. } \left| \frac{\mathbf{a}_i^\top \mathbf{z}}{\lambda} \right| \leq 1; s(-\mathbf{a}_i^\top \mathbf{z} + \mathbf{a}_i^\top \mathbf{w}^*) \leq 0 \right\} \\ &= \left\{ s \mid \forall \mathbf{z} \text{ s.t. } |\mathbf{a}_i^\top \mathbf{z}| \leq \lambda; s(\mathbf{a}_i^\top \mathbf{z}) \geq s(\mathbf{a}_i^\top \mathbf{w}^*) \right\} \end{aligned}$$

Since the optimality condition (4b) holds hence $-x_i^*$ should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g_i^*(-\mathbf{a}_i^\top \mathbf{w}^*)$ according to conditions given above. That is

$$-x_i^*(\mathbf{a}_i^\top \mathbf{z}) \leq -x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) \quad \forall \mathbf{z} \text{ s.t. } |\mathbf{a}_i^\top \mathbf{z}| \leq \lambda \quad (100)$$

$$x_i^*(\mathbf{a}_i^\top \mathbf{z}) \geq x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) \quad \forall \mathbf{z} \text{ s.t. } |\mathbf{a}_i^\top \mathbf{z}| \leq \lambda \quad (101)$$

$$\begin{aligned} \Rightarrow x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) &\leq \min_z (x_i^*(\mathbf{a}_i^\top \mathbf{z})) \quad \text{s.t. } |\mathbf{a}_i^\top \mathbf{z}| \leq \lambda \\ & \quad (102) \end{aligned}$$

Case 1: $x_i^* > 0$.

$$\begin{aligned} x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) &\leq \min_z (x_i^*(\mathbf{a}_i^\top \mathbf{z})) \quad \text{s.t. } |\mathbf{a}_i^\top \mathbf{z}| \leq \lambda \\ \Rightarrow x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) &\leq -\lambda x_i^* \\ \Rightarrow (\mathbf{a}_i^\top \mathbf{w}^*) &\leq -\lambda \\ \Rightarrow (\mathbf{a}_i^\top \mathbf{w}^*) &= -\lambda \end{aligned} \quad (103)$$

Equation (103) comes from the fact that $|\mathbf{a}_i^\top \mathbf{w}^*| \leq \lambda$

Case 2: $x_i^* < 0$.

$$\begin{aligned} x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) &\leq \min_z (x_i^*(\mathbf{a}_i^\top \mathbf{z})) \quad \text{s.t. } |\mathbf{a}_i^\top \mathbf{z}| \leq \lambda \\ \Rightarrow x_i^*(\mathbf{a}_i^\top \mathbf{w}^*) &\leq \lambda x_i^* \\ \Rightarrow (\mathbf{a}_i^\top \mathbf{w}^*) &\geq \lambda \\ \Rightarrow (\mathbf{a}_i^\top \mathbf{w}^*) &= \lambda \end{aligned} \quad (104)$$

Equation (103) comes from the fact that $|\mathbf{a}_i^\top \mathbf{w}^*| \leq \lambda$

Case 3: $x_i^* = 0$.

Since if we assume f as a continuous smooth function then $\mathbf{a}_i^\top \mathbf{w}^*$ is also continuous. Now if we consider arguments given for $x_i^* < 0$ and $x_i^* > 0$ we conclude that $|\mathbf{a}_i^\top \mathbf{w}^*| = \lambda$ in all of the above two cases. Since $x_i^* = 0$ is in the domain of the function (A), hence at $x_i^* = 0$, $\mathbf{a}_i^\top \mathbf{w}^*$ will lie in the open range of $-\lambda$ to λ . Which implies whenever $|\mathbf{a}_i^\top \mathbf{w}^*| < \lambda$, then $x_i^* = 0$

Another view on the proof can be derived from the optimality condition (4a).

The optimization problem (99) can be taken as partially separable problem and from the optimality condition (4a) kk

$$-\mathbf{a}_i^\top \mathbf{w}^* \in \partial g_i(x_i^*) \quad (105)$$

$$\partial g_i(x_i^*) \in \begin{cases} \lambda \frac{x_i^*}{|x_i^*|} & \text{if } x_i \neq 0 \\ [-\lambda, \lambda] & \text{if } x_i = 0 \end{cases} \quad (106)$$

From equations (114) and (115) we conclude that if

$$|\mathbf{a}_i^\top \mathbf{w}^*| < \lambda \Rightarrow x_i^* = 0$$

□

Proof of Theorem 16. From Equation (1a), we know that $\mathbf{w}^* \in \partial f(A\mathbf{x}^*)$. Hence from Lemma 23,

$$\begin{aligned} |\mathbf{a}_i^\top \mathbf{w}^*| &= |\mathbf{a}_i^\top (\mathbf{w}^* - \mathbf{w} + \mathbf{w})| \\ &\leq |\mathbf{a}_i^\top \mathbf{w}| + |\mathbf{a}_i^\top (\mathbf{w}^* - \mathbf{w})| \\ &\leq |\mathbf{a}_i^\top \mathbf{w}| + \|\mathbf{a}_i\|_2 \|\mathbf{w}^* - \mathbf{w}\|_2 \\ &\leq |\mathbf{a}_i^\top \mathbf{w}| + \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \end{aligned} \tag{107}$$

Eq. (107) comes from Corollary 2. Now using Lemma 23 and equation (107), we get

$$|\mathbf{a}_i^\top \nabla f(A\mathbf{x})| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \Rightarrow x_i^* = 0$$

□

Penalized Lasso. Screening in this case can be derived from the existing “gap safe” paper [Ndiaye et al., 2015]. For completeness we here show that the same result follows from our Theorem 16:

Corollary 24. Penalized Lasso Consider an optimization problem of the form:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

Then the screening rule is given by: $|\mathbf{a}_i^\top (A\mathbf{x} - \mathbf{b})| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{2G(\mathbf{x})} \Rightarrow x_i^* = 0$.

Proof of Corollary 24. By observing the cost function for penalized lasso it can be concluded that

$$f(A\mathbf{x}) = \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2, \quad \mathbf{w} = A\mathbf{x} - \mathbf{b}, \quad \text{and } L = 1$$

Now results from Theorem 16 can be directly applied here and hence the screening rule becomes

$$|\mathbf{a}_i^\top (A\mathbf{x} - \mathbf{b})| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{2G(\mathbf{x})} \Rightarrow x_i^* = 0.$$

□

This result is known in the literature [Ndiaye et al., 2015], and we recover it using our proposed general approach in this paper by using Theorem 16.

Also, by applying same trick as mentioned after the end of proof of Corollary 15, we can show that we can get rid of the factor 2 here also. Here also it is not hard to see that primal and dual ((A) and (B)) both are 1 strongly convex in the dual variable \mathbf{w} . Hence by the same argument as made in the proof of Corollary 15, we get that

$$G(\mathbf{w}) \geq \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

And the improved screening rule comes out to be

$$|\mathbf{a}_i^\top (A\mathbf{x} - \mathbf{b})| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{G(\mathbf{x})} \Rightarrow x_i^* = 0.$$

Logistic Regression with L_1 -regularization

Corollary 25. Logistic Regression with L_1 -norm Penalization. *The optimization problem for logistic regression with L_1 regularizer can be written in the form of:*

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^n \log(\exp([A\mathbf{x}]_i) + 1) + \lambda \|\mathbf{x}\|_1 \quad (108)$$

And screening rule for above problem can be written as :

$$\left| \mathbf{a}_i^\top \left(\frac{\exp(A\mathbf{x})}{\exp(A\mathbf{x}) + 1} \right) \right| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{2G(\mathbf{x})} \Rightarrow \mathbf{x}_i^* = 0$$

where $\left(\frac{\exp(A\mathbf{x})}{\exp(A\mathbf{x}) + 1} \right)$ is element wise vector whose i_{th} element is $\left(\frac{\exp([A\mathbf{x}]_i)}{\exp([A\mathbf{x}]_i) + 1} \right)$

Proof. By observation we know that in equation (108)

$$f(A\mathbf{x}) = \sum_{i=1}^n \log(\exp([A\mathbf{x}]_i) + 1) \text{ and } \mathbf{w} \text{ is elementwise vector of } w_i \text{ s.t } w_i = \frac{\exp([A\mathbf{x}]_i)}{\exp([A\mathbf{x}]_i) + 1}$$

According to [Smith et al., 2015, Lemma 5], we get that the function $f(A\mathbf{x})$ is 1-smooth. Hence $L = 1$
Now from theorem 16, we derive the screening rule for logistic regression with L_1 -regularization which is

$$\left| \mathbf{a}_i^\top \left(\frac{\exp(A\mathbf{x})}{\exp(A\mathbf{x}) + 1} \right) \right| < \lambda - \|\mathbf{a}_i\|_2 \sqrt{2G(\mathbf{x})} \Rightarrow \mathbf{x}_i^* = 0$$

where $\left(\frac{\exp(A\mathbf{x})}{\exp(A\mathbf{x}) + 1} \right)$ is element wise vector whose i_{th} element is $\left(\frac{\exp([A\mathbf{x}]_i)}{\exp([A\mathbf{x}]_i) + 1} \right)$ This result is also known in the literature in [Ndiaye et al., 2015] (or see also Wang et al. [2014] for a similar approach) and we recover it using our proposed general approach in this paper by using Theorem 16. □

Elastic-net regularized regression

Proof of Corollary 17.

$$\begin{aligned} & \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda_2 \|\mathbf{x}\|_2^2 + \lambda_1 \|\mathbf{x}\|_1 \\ &= \frac{1}{2} [\mathbf{x}^\top A^\top A \mathbf{x} - 2\mathbf{b}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{b}] + \lambda_2 \mathbf{x}^\top \mathbf{x} + \lambda_1 \|\mathbf{x}\|_1 \\ &= \frac{1}{2} [\mathbf{x}^\top (A^\top A + 2\lambda_2 I) \mathbf{x} - 2\mathbf{b}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{b}] + \lambda_1 \|\mathbf{x}\|_1 \end{aligned} \quad (109)$$

Now consider $A^\top A + 2\lambda_2 I = Q^\top Q$ and choose vector \mathbf{m} such that $A^\top \mathbf{b} = Q^\top \mathbf{m}$. Hence line (109) can be written as

$$\begin{aligned} & \frac{1}{2} [\mathbf{x}^\top (A^\top A + 2\lambda_2 I) \mathbf{x} - 2\mathbf{b}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{b}] + \lambda_1 \|\mathbf{x}\|_1 \\ &= \frac{1}{2} [\mathbf{x}^\top Q^\top Q \mathbf{x} - 2\mathbf{m}^\top Q \mathbf{x} + \mathbf{m}^\top \mathbf{m} - \mathbf{m}^\top \mathbf{m} + \mathbf{b}^\top \mathbf{b}] + \lambda_1 \|\mathbf{x}\|_1 \\ &= \frac{1}{2} \|Q\mathbf{x} - \mathbf{m}\|_2^2 + \frac{1}{2} [\mathbf{b}^\top \mathbf{b} - \mathbf{m}^\top \mathbf{m}] + \lambda_1 \|\mathbf{x}\|_1 \end{aligned}$$

Now the optimization problem (24) can be written as

$$\min_{\mathbf{x}} \frac{1}{2} \|Q\mathbf{x} - \mathbf{m}\|_2^2 + \lambda_1 \|\mathbf{x}\|_1 \quad (110)$$

Now results from Corollary 24 can be directly applied to (110).

From observation, we know that $f(Q\mathbf{x}) = \frac{1}{2} \|Q\mathbf{x} - \mathbf{m}\|_2^2$, $\mathbf{w} = Q\mathbf{x} - \mathbf{m}$, and $L = 1$
Simplification,

$$\begin{aligned} |\mathbf{q}_i^\top (Q\mathbf{x} - \mathbf{m})| &= |\mathbf{q}_i^\top Q\mathbf{x} - \mathbf{q}_i^\top \mathbf{m}| \\ &= |\mathbf{q}_i^\top Q\mathbf{x} - \mathbf{a}_i^\top \mathbf{b}| \\ &= |(\mathbf{a}_i^\top A + 2\lambda_2 \mathbf{e}_i^\top) \mathbf{x} - \mathbf{a}_i^\top \mathbf{b}| \end{aligned} \quad (111)$$

$$\begin{aligned}
 |q_i| \sqrt{2G(\mathbf{x})} &= \sqrt{\mathbf{a}_i^\top \mathbf{a}_i + 2\lambda_2} \sqrt{2G(\mathbf{x})} \\
 &= \sqrt{2(\mathbf{a}_i^\top \mathbf{a}_i + 2\lambda_2)G(\mathbf{x})}
 \end{aligned} \tag{112}$$

Now using results from Corollary 24, equations (111) and (112), we get screening rules for elastic norm regularization regression problem as:

$$|(\mathbf{a}_i^\top A + 2\lambda_2 \mathbf{e}_i^\top) \mathbf{x} - \mathbf{a}_i^\top \mathbf{b}| < \lambda_1 - \sqrt{2(\mathbf{a}_i^\top \mathbf{a}_i + 2\lambda_2)(G(\mathbf{x}))} \Rightarrow x_i^* = 0.$$

□

Lemma 26 (Conjugate of the Elastic Net Regularizer [Lemma 6 [Smith et al., 2015]]). *For $\alpha \in (0, 1]$, the elastic net function $g(\mathbf{x}) = \frac{1-\alpha}{2} \|\mathbf{x}\|_2^2 + \alpha \|\mathbf{x}\|_1$ is the convex conjugate of*

$$g^*(\mathbf{x}) = \sum_i \left[\frac{1}{2(1-\alpha)} ([|x_i| - \alpha]_+)^2 \right] = \sum_i g_i^*(x_i)$$

where $g_i(\beta_i) = [\frac{1-\alpha}{2} \beta_i^2 + \alpha |\beta_i|]$ and $[\cdot]_+$ is the positive part operator, $[s]_+ = s$ for $s > 0$, and zero otherwise. Furthermore, this g^* is smooth, i.e. has Lipschitz continuous gradient with constant $1/(1-\alpha)$.

Proof. The complete proof has been given in [Smith et al., 2015, Lemma 6] but we also provide proof here below. From the definition of convex conjugate function,

$$\begin{aligned}
 g^*(\mathbf{x}) &= \sup_{\boldsymbol{\beta}} [\mathbf{x}^\top \boldsymbol{\beta} - g(\boldsymbol{\beta})] \\
 &= \sup_{\boldsymbol{\beta}} \left[\mathbf{x}^\top \boldsymbol{\beta} - \left(\frac{1-\alpha}{2} \|\boldsymbol{\beta}\|_2^2 + \alpha \|\boldsymbol{\beta}\|_1 \right) \right] \\
 &= \sup_{\beta_i} \left[\sum_i x_i \beta_i - \left(\sum_i \left(\frac{1-\alpha}{2} \beta_i^2 + \alpha |\beta_i| \right) \right) \right] \quad \forall i \in [n] \\
 &= \sum_i \sup_{\beta_i} [x_i \beta_i - \left(\frac{1-\alpha}{2} \beta_i^2 + \alpha |\beta_i| \right)] \quad \forall i \in [n] \\
 &= \sum_i g_i^*(x_i), \text{ where } g_i(\beta_i) = \frac{1-\alpha}{2} \beta_i^2 + \alpha |\beta_i|
 \end{aligned}$$

Now,

$$g_i^*(x_i) = \sup_{\beta_i} [x_i \beta_i - \left(\frac{1-\alpha}{2} \beta_i^2 + \alpha |\beta_i| \right)]$$

Consider three cases now :

Case 1: $\beta > 0$.

$$\begin{aligned}
 g_i^*(x_i) &= \sup_{\beta_i} [x_i \beta_i - \left(\frac{1-\alpha}{2} \beta_i^2 + \alpha \beta_i \right)] \\
 &\Rightarrow \beta_i = \frac{(x_i - \alpha)}{(1-\alpha)} \text{ that also implies } x_i > \alpha \\
 \text{Hence, } g_i^*(x_i) &= \frac{(x_i - \alpha)^2}{2(1-\alpha)} \text{ whenever } x_i > \alpha
 \end{aligned}$$

Case 2: $\beta < 0$.

$$\begin{aligned}
 g_i^*(x_i) &= \sup_{\beta_i} [x_i \beta_i - \left(\frac{1-\alpha}{2} \beta_i^2 - \alpha \beta_i \right)] \\
 &\Rightarrow \beta_i = \frac{(x_i + \alpha)}{(1-\alpha)} \text{ that also implies } x_i < -\alpha \\
 \text{Hence, } g_i^*(x_i) &= \frac{(x_i + \alpha)^2}{2(1-\alpha)} \text{ whenever } x_i < -\alpha
 \end{aligned}$$

Case 3: $\beta = 0$.

$$g_i^*(x_i) = 0 \text{ that also implies } |x_i| \leq \alpha$$

Hence,

$$g_i^*(x_i) = \frac{1}{2(1-\alpha)} ([|x_i| - \alpha]_+)^2$$

From all of the above arguments, $g^*(\mathbf{x}) = \sum_i \left[\frac{1}{2(1-\alpha)} ([|x_i| - \alpha]_+)^2 \right] = \sum_i g_i^*(x_i)$ \square

Theorem 27. *If we consider the general elastic net formulation of the form*

$$\min_{\mathbf{x}} f(A\mathbf{x}) + \frac{(1-\alpha)}{2} \|\mathbf{x}\|_2^2 + \alpha \|\mathbf{x}\|_1 \quad (113)$$

If f is L -smooth, then the following screening rule holds for all $i \in [n]$:

$$|\mathbf{a}_i^\top \nabla f(A\mathbf{x})| < \alpha - \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \Rightarrow \mathbf{x}_i^* = 0$$

Proof. Since the optimization problem (113) comes under the partially separable framework and we can use the first order optimality condition (4a) as well as (4b) to derive screening rules for the problem.

By optimality condition (4b), we know that

$$x_i \in \partial g_i^*(-\mathbf{a}_i^\top \mathbf{w})$$

From lemma 26, $g_i^*(-\mathbf{a}_i^\top \mathbf{w}^*) = \frac{1}{2(1-\alpha)} ([| \mathbf{a}_i^\top \mathbf{w}^* | - \alpha]_+)^2$ and also $\partial g_i^*(-\mathbf{a}_i^\top \mathbf{w}) = 0$ whenever $| \mathbf{a}_i^\top \mathbf{w} | \leq \alpha$. Hence whenever $| \mathbf{a}_i^\top \mathbf{w} | \leq \alpha \Rightarrow x_i = 0$.

The same screening rule for elastic net regularized problem can be derived from the optimality condition (4a). The optimization problem (113) can be taken as partially separable problem and from the optimality condition (4a)

$$-\mathbf{a}_i^\top \mathbf{w}^* \in \partial g_i(x_i^*) \quad (114)$$

$$\partial g_i(x_i^*) \in \begin{cases} \alpha \frac{x_i^*}{|x_i^*|} + (1-\alpha)x_i & \text{if } x_i \neq 0 \\ [-\alpha, \alpha] & \text{if } x_i = 0 \end{cases} \quad (115)$$

Hence, whenever $| \mathbf{a}_i^\top \mathbf{w} | \leq \alpha \Rightarrow x_i = 0$.

The above arguments also show the significance of symmetry in our formulation as structure (A) and (B). This formulation provides our framework more flexibility to be used in larger class of problem.

Now,

$$\begin{aligned} | \mathbf{a}_i^\top \mathbf{w}^* | &= | \mathbf{a}_i^\top (\mathbf{w}^* - \mathbf{w} + \mathbf{w}) | \\ &\leq | \mathbf{a}_i^\top \mathbf{w} | + | \mathbf{a}_i^\top (\mathbf{w}^* - \mathbf{w}) | \\ &\leq | \mathbf{a}_i^\top \mathbf{w} | + \|\mathbf{a}_i\|_2 \|\mathbf{w}^* - \mathbf{w}\|_2 \\ &\leq | \mathbf{a}_i^\top \mathbf{w} | + \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \end{aligned} \quad (116)$$

Equation (116) comes directly from corollary 2. Hence finally we get the screening rules for general elastic net penalty problem which is very similar to screening for L_1 -penalized problems:

$$|\mathbf{a}_i^\top \nabla f(A\mathbf{x})| < \alpha - \|\mathbf{a}_i\|_2 \sqrt{2LG(\mathbf{x})} \Rightarrow \mathbf{x}_i^* = 0$$

Now the above mentioned rule can be made a bit tighter under some condition which is not very interesting to discuss \square

D.2 Screening for Structured Norms

Lemma 28. *If we use the same notation as mentioned in Section 5.3 to write a vector \mathbf{x} as a concatenation of smaller group vectors $\{\mathbf{x}_1 \cdots \mathbf{x}_G\}$ such that $\mathbf{x}^\top = [\mathbf{x}_1^\top, \mathbf{x}_2^\top \cdots \mathbf{x}_G^\top]$ and correspondingly the matrix A can be denoted as the concatenation of column groups $A = [A_1 \ A_2 \cdots A_G]$. Now if we consider an optimization problem of the form*

$$\arg \min_{\mathbf{x}} f(A\mathbf{x}) + \sum_{g=1}^G \sqrt{\rho_g} \|\mathbf{x}_g\|_2$$

At the optimal point \mathbf{x}^* and dual optimal points \mathbf{w}^* , we get rules according to the following equation:

$$\|A_g^\top \mathbf{w}^*\|_2 < \sqrt{\rho_g} \Rightarrow \mathbf{x}_g^* = 0$$

Proof. Dual of the problem is given by

$$\mathcal{O}_B(\mathbf{w}) = f^*(\mathbf{w}) + \sum_g \sqrt{\rho_g} \iota_{L_\infty} \left(\frac{\|A_g^\top \mathbf{w}\|_2}{\sqrt{\rho_g}} \right) \quad (117)$$

Hence for the indicator function g_g^* by Lemma 22

$$\begin{aligned} \partial g_g^*(-A_g^\top \mathbf{w}^*) &= \left\{ \mathbf{s} \mid \forall \mathbf{z} \text{ s.t. } \left\| \frac{A_g^\top \mathbf{z}}{\sqrt{\rho_g}} \right\|_2 \leq 1; \mathbf{s}^\top (-A_g^\top \mathbf{z} + A_g^\top \mathbf{w}^*) \leq 0 \right\} \\ &= \left\{ \mathbf{s} \mid \forall \mathbf{z} \text{ s.t. } \|A_g^\top \mathbf{z}\|_2 \leq \sqrt{\rho_g}; \mathbf{s}^\top (A_g^\top \mathbf{z}) \geq \mathbf{s}^\top (A_g^\top \mathbf{w}^*) \right\} \end{aligned}$$

Now, by the optimality condition (4b) $\mathbf{x}_g \in \partial g_g^*(-A_g^\top \mathbf{w}^*)$, and since this holds, hence xv_g^* should satisfy the required constrained which is needed to be in the set of subgradients of $\partial g^*(-A_g^\top \mathbf{w}^*)$ according to conditions given above. Hence,

$$\begin{aligned} -\mathbf{x}_g^{*\top} (A_g^\top \mathbf{z}) &\leq -\mathbf{x}_g^{*\top} (A_g^\top \mathbf{w}^*) \quad \forall \mathbf{z} \text{ s.t. } \|A_g^\top \mathbf{z}\|_2 \leq \sqrt{\rho_g} \\ \Rightarrow \mathbf{x}_g^{*\top} (A_g^\top \mathbf{z}) &\geq \mathbf{x}_g^{*\top} (A_g^\top \mathbf{w}^*) \quad \forall \mathbf{z} \text{ s.t. } \|A_g^\top \mathbf{z}\|_2 \leq \sqrt{\rho_g} \\ \Rightarrow \mathbf{x}_g^{*\top} (A_g^\top \mathbf{w}^*) &\leq \min_z \mathbf{x}_g^{*\top} (A_g^\top \mathbf{z}) \quad \text{s.t. } \|A_g^\top \mathbf{z}\|_2 \leq \sqrt{\rho_g} \\ \Rightarrow \mathbf{x}_g^{*\top} (A_g^\top \mathbf{w}^*) &\leq \min_z \|\mathbf{x}_g\|_2 \|A_g^\top \mathbf{z}\|_2 \quad \text{s.t. } \|A_g^\top \mathbf{z}\|_2 \leq \sqrt{\rho_g} \\ \Rightarrow \mathbf{x}_g^{*\top} (A_g^\top \mathbf{w}^*) &\leq -\|\mathbf{x}_g^*\|_2 \sqrt{\rho_g} \\ \Rightarrow \|A_g^\top \mathbf{w}^*\|_2 &= \sqrt{\rho_g} \end{aligned} \quad (118)$$

Equation (118) comes from the cauchy inequality and true $\forall \mathbf{x}_g^* : \mathbf{x}_g^* \neq 0$. Whenever $\|A_g^\top \mathbf{w}^*\|_2 < \sqrt{\rho_g}$ then $\mathbf{x}_g^* = 0$

Another view on the screening of above optimization problem can be seen from the optimality condition (4a). The optimization problem in Lemma 28 can be taken as partially separable problem and from the optimality condition (4a)

$$-A_g^\top \mathbf{w}^* \in \partial g(\mathbf{x}_g^*) \quad (119)$$

$$\partial g(\mathbf{x}_g^*) \in \begin{cases} \sqrt{\rho_g} \frac{\mathbf{x}_g}{\|\mathbf{x}_g\|_2} & \text{if } \mathbf{x}_g \neq 0 \\ \mathcal{B}_2 & \text{if } \mathbf{x}_g = 0 \text{ and } \mathcal{B}_2 \text{ is norm ball of radius } \sqrt{\rho_g} \end{cases} \quad (120)$$

From Equations (119) and (120), we conclude that if

$$\|A_g^\top \mathbf{w}^*\|_2 < \sqrt{\rho_g} \Rightarrow \mathbf{x}_g^* = 0$$

□

Proof of Theorem 18. From Equation (1a), we know that $\mathbf{w} \in \nabla f(A\mathbf{x})$. Now

$$\begin{aligned}
 \|A_g^\top \mathbf{w}^*\|_2 &= \|A_g^\top (\mathbf{w} + \mathbf{w}^* - \mathbf{w})\|_2 \leq \|A_g^\top \mathbf{w}\|_2 + \|A_g^\top (\mathbf{w}^* - \mathbf{w})\|_2 \\
 &= \|A_g^\top \mathbf{w}\|_2 + \sqrt{\text{tr}((A_g^\top (\mathbf{w}^* - \mathbf{w}))((\mathbf{w}^* - \mathbf{w})^\top)A_g)} \\
 &\leq \|A_g^\top \mathbf{w}\|_2 + \sqrt{\text{tr}((\mathbf{w}^* - \mathbf{w})^\top (\mathbf{w}^* - \mathbf{w}))} \sqrt{\text{tr}(A_g^\top A_g)} \\
 &= \|A_g^\top \mathbf{w}\|_2 + \|\mathbf{w}^* - \mathbf{w}\|_2 \|A_g\|_{\text{Fro}}
 \end{aligned} \tag{121}$$

Using Corollary 2 with Equation (121), we get

$$\|A_g^\top \mathbf{w}^*\|_2 \leq \|A_g^\top \nabla f(A\mathbf{x})\|_2 + \sqrt{2LG(\mathbf{x})} \|A_g\|_{\text{Fro}}$$

Hence using previous Lemma 28,

$$\|A_g^\top \nabla f(A\mathbf{x})\|_2 + \sqrt{2LG(\mathbf{x})} \|A_g\|_{\text{Fro}} < \sqrt{\rho_g} \Rightarrow \mathbf{x}_g^* = 0$$

□

Proof of Corollary 19. This is an explicit case of the optimization problem mentioned in Lemma 28. By observation we know that,

$$f(A\mathbf{x}) = \frac{1}{2} \|A\mathbf{x} - b\|^2, \quad \mathbf{w} = A\mathbf{x} - b \quad \text{and} \quad L = 1$$

Now applying the findings of Theorem 18, we get

$$\|A_g^\top (A\mathbf{x} - b)\|_2 + \sqrt{2G(\mathbf{x})} \|A_g\|_{\text{Fro}} < \lambda \sqrt{\rho_g} \Rightarrow \mathbf{x}_g^* = 0$$

□

In Lemma 29 mentioned below, we show that the structured norm setting of [Ndiaye et al., 2015] can be derived from our more general (A) and (B) structure.

Lemma 29. *Sparse Multi-Task and Multi Class Model [Ndiaye et al., 2015] - If we consider general problem of the form*

$$\min_{X \in \mathbb{R}^{p \times q}} \sum_{i=1}^n f_i(\mathbf{a}_i^\top X) + \lambda \Omega(X) \tag{122}$$

where the regularization function $\Omega : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}_+$ is such that $\Omega(X) = \sum_{g=1}^p \|\mathbf{x}_g\|_2$ and $X = [\mathbf{x}_1, \mathbf{x}_2 \cdots \mathbf{x}_G]$. We write $W = [\mathbf{w}_1, \mathbf{w}_2 \cdots \mathbf{w}_G]$ for variable of the dual problem. Then the screening rule becomes

$$\|\mathbf{a}^{(g)\top} W\|_2 < \lambda - \|\mathbf{a}^{(g)}\|_2 \|W - W^*\|_2 \Rightarrow \mathbf{x}_g^* = 0$$

Here $\mathbf{a}^{(g)}$ is the vector of the g^{th} element group of each vector \mathbf{a}_i .

Proof. Equations pair (A) and (B) can be used interchangeably by replacing primal with dual and f with g . Hence the partial separable primal-dual pair (SA) and (SB) can also be used interchangeably. By comparing Equation (122) with (SA) and (SB), we observe that separable function $\sum_{i=1}^n f_i(\mathbf{a}_i^\top X)$ takes the place of separable g^* in (SB) and $\lambda \Omega(X)$ takes the place of f^* . Hence we apply the optimality condition (1b) to get (with exchanged primal dual variable)

$$AW^* \in \partial \lambda \Omega(X^*)$$

Hence if,

$$\|\mathbf{a}^{(g)\top} W^*\|_2 < \lambda \Rightarrow \mathbf{x}_g = 0 \tag{123}$$

Now,

$$\begin{aligned}
 \|\mathbf{a}^{(g)\top} W^*\|_2 &= \|\mathbf{a}^{(g)\top} (W^* - W + W)\|_2 \\
 &\leq \|\mathbf{a}^{(g)\top} W\|_2 + \|\mathbf{a}^{(g)\top} (W^* - W)\|_2 \\
 &\leq \|\mathbf{a}^{(g)\top} W\|_2 + \|\mathbf{a}^{(g)}\|_2 \|W^* - W\|_2
 \end{aligned} \tag{124}$$

Using equations (123) and (124), the screening rule comes out to be

$$\|\mathbf{a}^{(g)\top} W\|_2 < \lambda - \|\mathbf{a}^{(g)}\|_2 \|W - W^*\|_2 \Rightarrow \mathbf{x}_g^* = 0$$

□

Corollary 30. *If for all $i \in [n]$, f_i is L -Lipschitz gradient then screening rule for equation (122) is*

$$\|\mathbf{a}^{(g)\top} W\|_2 < \lambda - \|\mathbf{a}^{(g)}\|_2 \sqrt{2L G(X)} \Rightarrow \mathbf{x}_g^* = 0$$

Proof. Using Lemma 29 and Corollary 2, we get the desired expression. □

D.3 Connection with Sphere Test Method

The general idea behind the sphere test method Xiang et al. [2014] is to consider the maximum value of desired function in a spherical region which contains the optimal dual variable. In context of our general framework (A) and (B), we obtain this case when considering an ℓ_1 penalty or ℓ_2/ℓ_1 penalty. That means g is a norm and hence from Lemma 20, g^* becomes the indicator function of the dual norm ball of $A^\top \mathbf{w}$. The dual norm function for ℓ_1 norm is of the form $\max_i |\mathbf{a}_i^\top \mathbf{w}|$ and for ℓ_2/ℓ_1 norm, it is $\max_g \|A_g^\top \mathbf{w}\|$. Hence, we try to find maximum value of the function of the forms $\max_{\theta \in \mathcal{S}(\mathbf{q}, r)} \mathbf{a}_i^\top \theta$ where $\mathcal{S}(\mathbf{q}, r) = \{z : \|z - \mathbf{q}\|_2 \leq r\}$ the ball \mathcal{S} also contains the optimal dual point \mathbf{w}^* . If the maximum value of $\mathbf{a}_i^\top \theta$ is less than some particular value for all the θ in the ball hence $\mathbf{a}_i^\top \mathbf{w}$ will also be less than that particular value and that is the main reason we try to find maximum of $\mathbf{a}_i^\top \theta$ over the ball \mathcal{S} .

$$\begin{aligned} \max_{\theta \in \mathcal{S}(\mathbf{q}, r)} \mathbf{a}_i^\top \theta &= \mathbf{a}_i^\top (\theta - \mathbf{q} + \mathbf{q}) = \mathbf{a}_i^\top (\theta - \mathbf{q}) + \mathbf{a}_i^\top \mathbf{q} \\ &\leq \|\mathbf{a}_i\|_2 \|\theta - \mathbf{q}\| + \mathbf{a}_i^\top \mathbf{q} \leq r \|\mathbf{a}_i\|_2 + \mathbf{a}_i^\top \mathbf{q} \end{aligned}$$

Similar arguments can be given in the ℓ_2/ℓ_1 -norm case. A variety of existing screening test for lasso and group lasso are of this flavor of sphere tests. The difference between these approaches mainly lie in the way of choosing the center and bounding the radius of the sphere, such that the optimal dual variables lie inside the sphere. Our method can be seen as a general framework for such a sphere test based screening with dynamic screening rules. Our method can be interpreted as a sphere test with the current iterate of the dual variable \mathbf{w} as a center of the ball, and we obtain the bound on the radius in terms of duality gap function.