Solving large-scale L1-regularized SVMs and cousins: the surprising effectiveness of column and constraint generation

Antoine Dedieu* Rahul Mazumder[†]

December, 2018

Abstract

The linear Support Vector Machine (SVM) is one of the most popular binary classification techniques in machine learning. Motivated by applications in modern high dimensional statistics, we consider penalized SVM problems involving the minimization of a hinge-loss function with a convex sparsity-inducing regularizer such as: the L1-norm on the coefficients, its grouped generalization and the sorted L1-penalty (aka Slope). Each problem can be expressed as a Linear Program (LP) and is computationally challenging when the number of features and/or samples is large – the current state of algorithms for these problems is rather nascent when compared to the usual L2-regularized linear SVM. To this end, we propose new computational algorithms for these LPs by bringing together techniques from (a) classical column (and constraint) generation methods and (b) first order methods for non-smooth convex optimization — techniques that are rarely used together for solving large scale LPs. These components have their respective strengths; and while they are found to be useful as separate entities, they have not been used together in the context of solving large scale LPs such as the ones studied herein. Our approach complements the strengths of (a) and (b) — leading to a scheme that seems to outperform commercial solvers as well as specialized implementations for these problems by orders of magnitude. We present numerical results on a series of real and synthetic datasets demonstrating the surprising effectiveness of classic column/constraint generation methods in the context of challenging LP-based machine learning tasks.

1 Introduction

The linear Support Vector Machine (SVM) [34, 17] is an extremely popular and useful tool for binary classification. Given training data $(\mathbf{x}_i, y_i)_{i=1}^n$ with feature vector $\mathbf{x}_i \in \mathbb{R}^p$ and label $y_i \in \{-1, 1\}$, the task is to learn a linear classifier of the form $\operatorname{sign}(\mathbf{x}^T\boldsymbol{\beta} + \beta_0)$ where, $\beta_0 \in \mathbb{R}$ is the offset. The popular L2-regularized linear SVM (L2-SVM) considers the minimization problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \beta_0 \in \mathbb{R}} \sum_{i=1}^n \left(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \right)_+ + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2$$
 (1)

where, $(a)_{+} := \max\{a, 0\}$ is often noted as the hinge-loss function. Several algorithms have been proposed to efficiently solve Problem (1). Popular approaches include stochastic subgradient methods on the primal form [8, 29], coordinate descent methods on a dual [18] and cutting plane

^{*[}adedieu@mit.edu]. Operations Research Center, MIT

^{† [}rahulmaz@mit.edu]. MIT Sloan School of Management and Operations Research Center, MIT

algorithms [20, 14]. The L2-SVM estimator leads to a dense estimate for β — towards this end, the L1 penalty [11, 17] is often used as a convex surrogate to encourage sparsity (i.e., few nonzeros) in the coefficients. This leads to the L1-SVM problem:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \beta_0 \in \mathbb{R}} \sum_{i=1}^n \left(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \right)_+ + \lambda \|\boldsymbol{\beta}\|_1, \tag{2}$$

which can be written as a Linear Program (LP). The regularization parameter $\lambda \geq 0$ controls the degree of shrinkage on β . Off-the-shelf solvers, including commercial LP solvers (eg, Gurobi, Cplex) work very well for small/moderate sized problems, but have difficulties to solve Problem (2) when n and/or p is large (around ten thousand or so). Leading specialized solvers for Problem (2) include: a homotopy based method to compute the entire (piecewise linear) regularization path in β [16]; an Alternating Direction Method of Multipliers (ADMM) [2] based method. Recently [26] proposes a parametric simplex approach to solve Problem (2), which seems to be the current state-of-the-art. However, as our experiments suggest — the run times of these algorithms increase with problem-size, and the computations become prohibitively expensive (especially when compared to the algorithms we propose herein) as soon as $p \approx 5000$ or more — thereby seriously limiting the use of Problem (2) in tasks that arise in practice. In this paper, we address this shortcoming by bringing to bear somewhat underutilized classic operations research tools such as cutting planes (e.g., column/constraint generation), and popular modern first order optimization techniques—thereby presenting a new approach to solve Problem (2) and its cousins, introduced below.

In several applications, sparsity is structured — the coefficient indices are naturally found to occur in groups that are known a-priori and it is desirable to select (or set to zero) a whole group together as a "unit". In this context, a group version of the usual L1 norm is often used to improve the performance and interpretability of the model [36, 19]. We consider the popular L1/ L_{∞} penalty [1] leading to the Group-SVM Problem:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \ \beta_0 \in \mathbb{R}} \sum_{i=1}^n \left(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \right)_+ + \lambda \sum_{g=1}^G \|\boldsymbol{\beta}_g\|_{\infty}$$
 (3)

where, g = 1, ..., G denotes a group index (the groups are disjoint), β_g denotes the subvector of coefficients belonging to group g and $\beta = (\beta_1, ..., \beta_G)$. Problem (3) can be expressed as an LP and our approach applies to this problem as well.

The third problem we study in this paper is of a rather different flavor and is inspired by the sorted L1-penalty aka the Slope norm [7, 5], popularly used in the context of penalized least squares problems for its useful statistical properties. For a (a-priori specified) sequence $\lambda_1 \geq \ldots \geq \lambda_p \geq 0$,

the Slope-SVM problem is given by:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \ \beta_0 \in \mathbb{R}} \sum_{i=1}^n \left(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \right)_+ + \sum_{j=1}^p \lambda_j |\beta_{(j)}|, \tag{4}$$

where $|\beta_{(1)}| \geq \ldots \geq |\beta_{(p)}|$ are the ordered values of $|\beta_i|, i = 1, \ldots, p$. We show in Section 3 that Problem (4) can be expressed as an LP with O(n+p) variables and an exponential number (in p) of constraints, consequently posing challenges in optimization — therefore, specialized algorithms are called for. Using standard reformulation methods [9] (see Section A.2), Problem (4) can be modeled (e.g., using CVXPY) and solved (e.g., using a commercial solver like Gurobi) for small-sized problems. However, the computations become expensive when λ_i s are distinct (e.g., as in [7]) — for these cases, CVXPY can handle problems up to n = 100, p = 200 whereas, our approach can solve problems with $p \approx 50,000$ within a minute.

First order methods [23] have enjoyed great success in solving large scale structured convex optimization problems arising in machine learning applications. Many of these methods (such as proximal gradient and its accelerated variants) are appealing for minimization of smooth functions and also problems of the composite form [25], wherein they enjoy a convergence rate of $O(1/\sqrt{\epsilon})$ to obtain an ϵ -accurate solution. These algorithms however, do not directly apply to the nonsmooth SVM problems ((2),(3),(4)) discussed above. Nesterov's smoothing method [24] (which replaces the hinge-loss with a smooth approximation) can be used to obtain algorithms with a convergence rate of $O(1/\epsilon)$ —a method that we explore in Section 4. While this procedure (with additional heuristics based on [30]) lead to low accuracy solutions relatively fast; in our experience, the basic version of this algorithm takes a long time to obtain a solution with higher accuracy when nand/or p are large. Indeed a similar story applies to first order methods based on [4] and [10, 2] - their run times increase as the problem sizes become large. Since a main purpose of this paper is to demonstrate the power of classical cutting plane techniques (e.g., column and constraint generation), we focus on deterministic algorithms, as opposed to stochastic algorithms¹. We have observed empirically that the L1-SVM problem can be solved quite efficiently with commercial LP solvers (e.g., Gurobi) for instances with $n \approx 100$ and $p \approx 10^4$. Curiously, this is faster than standard or existing implementations of first order methods (e.g., based on ADMM [2] or Nesterov's smoothing method) to obtain solutions of a similar accuracy. However, these run times increase steeply when the problem sizes become large in n and/or p, thereby necessitating new algorithmic approaches such as those proposed herein.

What this paper is about: In this paper, we propose an efficient algorithmic framework for L1-SVM (and its relatives) leveraging the efficiency of excellent off-the-shelf LP solvers and in addition,

¹The cutting plane algorithms used here are deterministic methods and hence allows for a fair comparison.

exploiting the structural properties of solutions to these problems. For example, large values of λ will encourage an optimal solution to Problem (2), $\hat{\beta}$ (say), to be sparse. This sparsity will be critical to solve Problem (2) when $p \gg n$ —we anticipate to solve Problem (2) without having to create an LP model with all p variables. To this end, we will use column generation methods, a classical method in mathematical optimization/operations research with origins dating back to at least 1958 [13, 12, 6] — these techniques are commonly used to solve large scale LPs where many of the variables are anticipated to be zero. These methods however, to our knowledge, have received limited attention in the context of L1-SVM problems. Evidence presented herein suggests that these are very effective algorithms—we advocate that they be used more frequently to solve machine learning tasks based on LPs, even beyond the ones studied here.

We also consider another special structural aspect of a solution to Problem (2) when n is large (and p is not too large). Here, at an optimal solution, many of the points y_i 's will be correctly classified — that is, $\alpha_i := 1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) < 0$ and hence $\tilde{\alpha}_i := (\alpha_i)_+$ will be zero for many indices $i = 1, \ldots, n$. We leverage this sparsity in $\tilde{\alpha}_i$'s to develop efficient algorithms for Problem (2), using constraint generation [13, 12] methods. This observation allows us to solve Problem (2) without explicitly creating an LP model with as many samples.

To summarize, there are two characteristics special to an optimal solution of Problem (2): (a) sparsity in the SVM coefficients, i.e., β and/or (b) sparsity in $\tilde{\alpha}_i$'s. Column generation can be used to handle (a); constraint generation can be used to address (b) — in problems where both n, p are large, we propose to combine both column and constraint generation. To our knowledge, while column generation and constraint generation are used separately in the context of solving large scale LPs, using them together, especially in the context of the current application, is novel. For solving these (usually small) subproblems, we rely on powerful LP solvers (e.g., simplex based algorithms of Gurobi) which also have excellent warm-starting capabilities.

The cutting plane methods mentioned above, are found to benefit from good initializations. To this end, we use first order optimization methods to get approximate solutions with low computational cost. These solutions serve as decent initializations and are subsequently improved to deliver optimal solutions as a part of our column and/or constraint generation framework. Our approach also applies to the Group-SVM Problem (3). We extend our approach in a non-trivial manner to address the Slope-SVM problem (4). To our knowledge, this is the first paper that discusses the notion of bringing together techniques from first order methods in convex optimization and cutting plane algorithms for solving large scale LPs, in the context of solving a problem of key importance in machine learning. Implementation of our methods can be found at:

https://github.com/antoine-dedieu/cutting_planes_11_SVM_and_cousins.

Organization of paper: The rest of this paper is organized as follows. Section 2 discusses col-

umn/constraint generation methods for the L1-SVM and Group-SVM problems. Section 3 discusses its generalization to Slope-SVM. Section 4 discusses how first order methods can be used to get approximate solutions for these problems. Section 5 presents numerical results.

Notation: For an integer a we use [a] to denote $\{1, 2, ..., a\}$. The ith entry of a vector \mathbf{u} is denoted by u_i . For a set \mathcal{A} , we use the notation $|\mathcal{A}|$ to denote its size. For a positive semidefinite matrix \mathbf{A} , we denote its largest eigenvalue by $\sigma_{\max}(\mathbf{A})$.

2 Cutting plane algorithms for L1-SVM and its group extension

2.1 Primal and dual formulations of L1-SVM

We present an LP formulation for Problem (2):

Above, the positive and negative parts of β_i are denoted as $\beta_i^+ = \max\{\beta_i, 0\}$ and $\beta_i^- = \max\{-\beta_i, 0\}$ respectively, and ξ_i 's are auxiliary continuous variables corresponding to the hinge-loss function. The feasible set of Problem (5) is nonempty. A dual [6] of (5) is the following LP:

(Dual-L1-SVM):
$$\max_{\boldsymbol{\pi} \in \mathbb{R}^n} \qquad \sum_{i=1}^n \pi_i$$
s.t.
$$-\lambda \leq \sum_{i=1}^n y_i x_{ij} \pi_i \leq \lambda \qquad j \in [p]$$

$$\mathbf{y}^T \boldsymbol{\pi} = 0$$

$$0 \leq \pi_i \leq 1 \qquad i \in [n].$$
(6)

For Problems (5) and (6), standard complementary slackness conditions lead to:

$$(1 - \pi_i)\xi_i = 0, \qquad \pi_i \left(\xi_i + y_i \mathbf{x}_i^T \boldsymbol{\beta} + y_i \beta_0 - 1 \right) = 0 \qquad i \in [n].$$
 (7)

Let $(\beta^*(\lambda), \beta_0^*(\lambda))$ and $\pi^*(\lambda)$ denote optimal solutions for Problems (5) and (6). In what follows, for notational convenience, we will drop the dependence (of an optimal solution) on λ when there is no confusion. In standard SVM terminology [17], the vectors on the correct side of the margin satisfy $\xi_i = 0$ and $\pi_i = 0$, the ones lying on the margin satisfy $\xi_i = 0$ and $0 < \pi_i < 1$, and those on

the wrong side satisfy $\xi_i > 0$ and $\pi_i = 1$. Samples lying on the margin or on the wrong side of the margin are the *support vectors*: they fully define the maximal separating hyperplane.

2.2 Column generation for L1-SVM

Section A.1 reviews column and constraint generation for a general LP — here we discuss how column generation might be applied to solve the L1-SVM Problem (2) for a given λ with n a few hundred and p up to a million or so². Given a set of candidate features (cf Section 2.2.1) $\mathcal{J} \subset \{1, \ldots, p\}$, we form the restricted columns L1-SVM problem as

$$\begin{bmatrix}
\mathcal{M}_{\ell_{1}}([n], \mathcal{J}) \end{bmatrix} \quad \min_{\substack{\boldsymbol{\xi} \in \mathbb{R}^{n}, \beta_{0} \in \mathbb{R} \\ \boldsymbol{\beta}^{+}, \ \boldsymbol{\beta}^{-} \in \mathbb{R}^{|\mathcal{J}|}}$$
s.t.
$$\sum_{i=1}^{n} \xi_{i} + \lambda \sum_{j \in \mathcal{J}} \beta_{j}^{+} + \lambda \sum_{j \in \mathcal{J}} \beta_{j}^{-}$$
s.t.
$$\xi_{i} + \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{+} - \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{-} + y_{i} \beta_{0} \geq 1 \qquad i \in [n]$$

$$\boldsymbol{\xi} \geq 0, \ \boldsymbol{\beta}^{+} \geq 0, \ \boldsymbol{\beta}^{-} \geq 0.$$
(8)

Similar to Problem (6), we can get a dual of the restricted problem $\mathcal{M}_{\ell_1}([n], \mathcal{J})$ — let $\pi^* \in \mathbb{R}^n$ be an optimal solution to this restricted dual. For variable β_j^+ , we denote its corresponding reduced cost (cf Section A.1) by $\bar{\beta}_j^+$. A similar notation is used for β_j^- . For every pair of variables β_j^+ , β_j^- , $j \notin \mathcal{J}$, the minimum of their reduced costs (cf Section A.1) is

$$\min\left\{\bar{\beta}_j^+, \bar{\beta}_j^-\right\} = \lambda - \left|\sum_{i=1}^n y_i x_{ij} \pi_i^*\right|. \tag{9}$$

For a tolerance level $\epsilon > 0$ (we use $\epsilon = 10^{-2}$ in experiments), we update \mathcal{J} by adding all columns corresponding to pairs (β_j^+, β_j^-) such that the minimum of their reduced costs is lower than $-\epsilon$. We continue till no further column can be added. We summarize the algorithm below.

Algorithm 1: Column generation for L1-SVM

Input: X, y, regularization parameter λ , a convergence threshold $\epsilon > 0$, a set of columns \mathcal{J} . Output: A near-optimal solution β^* for the L1-SVM Problem (2).

- 1. Repeat Steps 2 to 3 until \mathcal{J} stabilizes.
- 2. Solve the problem $\mathcal{M}_{\ell_1}([n], \mathcal{J})$ (cf Problem (8)).
- 3. Form the set \mathcal{J}^{ϵ} of columns in $\{1,\ldots,p\}\setminus\mathcal{J}$ with reduced cost lower than $-\epsilon$. Update $\mathcal{J}\leftarrow\mathcal{J}\cup\mathcal{J}^{\epsilon}$; and go to Step 2 (using LP warm-starting).

 $^{^{2}}$ When p is a few thousand, we observe that Gurobi LP solvers work quite well and hence the cutting plane techniques discussed are not critical.

2.2.1 Initializing column generation with a candidate set of columns

In practice, Algorithm 1 is found to benefit from a good initial choice for \mathcal{J} . To obtain a reasonable estimate of \mathcal{J} with low computational cost, we found the following schemes to be useful:

- (i) (Correlation screening) A simple approach based on correlation screening selects \mathcal{J} as a subset (of size close to n) of variables with highest absolute inner product³ with \mathbf{y} . Computational experiments are reported in Section 5.1.1.
- (ii) (Regularization path) Here we compute a path of solutions to L1-SVM (with column generation) for a decreasing sequence of λ values (e.g., $\lambda \in \{\lambda_0, \dots, \lambda_M\}$) with the smallest one set to the current value of interest. This method which is discussed in Section 2.2.2 can also be used to compute a regularization path for the L1-SVM problem via warm-start continuation.
 - (iii) (First order methods) Section 4 discusses first order methods to obtain \mathcal{J} .

2.2.2 Computing a regularization path with column generation

Note that the subgradient condition of optimality for the L1-SVM Problem (2) is given by:

$$\lambda \operatorname{sign}(\beta_j^*) = \sum_{i=1}^n y_i x_{ij} \pi_i^*$$

where, sign(u) denotes a subgradient of $u \mapsto |u|$. When λ is larger than $\lambda_{\max} = \max_{j \in [p]} \sum_{i=1}^{n} |x_{ij}|$, an optimal solution to Problem (2) is zero: $\boldsymbol{\beta}^*(\lambda) = \mathbf{0}$.

Let \mathcal{I}_+ , \mathcal{I}_- denote the sample indices corresponding to the classes with labels +1 and -1 (respectively); and let N_+ , N_- denote their respective sizes. If $N_+ \geq N_-$, then for $\lambda \geq \lambda_{\max}$ a solution to Problem (6) is $\pi_i(\lambda) = N_-/N_+$, $\forall i \in \mathcal{I}_+$ and $\pi_i(\lambda) = 1$, $\forall i \in \mathcal{I}_-$. For $\lambda = \lambda_{\max}$ using (9), the minimum of the reduced costs of the variables β_j^+ and β_j^- is

$$\min\left\{\bar{\beta}_{j}^{+}(\lambda_{\max}), \bar{\beta}_{j}^{-}(\lambda_{\max})\right\} = \lambda_{\max} - \left|\frac{N_{-}}{N_{+}} \sum_{i \in \mathcal{I}_{+}} y_{i} x_{ij} + \sum_{i \in \mathcal{I}_{-}} y_{i} x_{ij}\right|. \tag{10}$$

When $\lambda = \lambda_1$ is slightly smaller than λ_{max} , (10) suggests to run a column generation algorithm by selecting \mathcal{J} as a small subset of variables that minimize the right-hand side of (10). Once we obtain a solution to Problem (2) at λ_1 , we can compute a solution for a smaller value of λ by using the warm-start capabilities of a simplex-based LP solver, along with column generation. This can also be used to compute an entire regularization path as summarized below.

³When the features are standardized to have zero mean with unit L2-norm, this is equivalent to sorting the absolute correlations: We thus use the phrase 'correlation screening' with a slight abuse of terminology.

Algorithm 2: Regularization path algorithm for L1-SVM

Input: \mathbf{X} , \mathbf{y} , convergence tolerance ϵ , a grid of decreasing λ values: $\{\lambda_0 = \lambda_{\max}, \dots, \lambda_M = \lambda\}$, a small integer j_0 .

Output: A near-optimal regularization path $\{\beta^*(\lambda_0), \dots, \beta^*(\lambda_M)\}$ for the L1-SVM Problem (2).

- 1. Let $\beta^*(\lambda_0) = \mathbf{0}$ and assign $\mathcal{J}(\lambda_0)$ to the j_0 variables minimizing the rhs of (10).
- 2. For $\ell \in \{1, ..., M\}$ initialize $\mathcal{J}(\lambda_{\ell}) \leftarrow \mathcal{J}(\lambda_{\ell-1})$, $\beta^*(\lambda_{\ell}) \leftarrow \beta^*(\lambda_{\ell-1})$. Run the column generation algorithm to obtain the new estimate $\beta^*(\lambda_{\ell})$ with $\mathcal{J}(\lambda_{\ell})$ denoting the corresponding set of columns.

2.3 Column and constraint generation for L1-SVM

We now consider the case where n is large (e.g., hundreds of thousands). If p is small compared to n, for the L1-SVM Problem (2), a separating hyperplane (corresponding to a solution of the problem) can be described by a small number of samples. We plan to use constraint generation ideas (Section A.1) to exploit this characteristic. We first present the case where n is large but p is small (Section 2.3.1) and then discuss the case where both n, p are large (Section 2.3.2).

2.3.1 Constraint generation for large n and small p

Given $\lambda > 0$ and $\mathcal{I} \subset \{1, \dots, n\}$, we define the restricted constraints version of the L1-SVM problem, by using a subset (indexed by \mathcal{I}) of the constraints in Problem (5)

$$\underbrace{\mathcal{M}_{\ell_{1}}(\mathcal{I},[p])}_{\boldsymbol{\xi} \in \mathbb{R}^{|\mathcal{I}|},\beta_{0} \in \mathbb{R}} \qquad \sum_{i \in \mathcal{I}} \xi_{i} + \lambda \sum_{j=1}^{p} \beta_{j}^{+} + \lambda \sum_{j=1}^{p} \beta_{j}^{-}$$

$$s.t. \qquad \xi_{i} + \sum_{j=1}^{p} y_{i}x_{ij}\beta_{j}^{+} - \sum_{j=1}^{p} y_{i}x_{ij}\beta_{j}^{-} + y_{i}\beta_{0} \ge 1 \qquad i \in \mathcal{I}$$

$$\boldsymbol{\xi} \ge 0, \ \boldsymbol{\beta}^{+} \ge 0, \ \boldsymbol{\beta}^{-} \ge 0.$$
(11)

A dual of Problem (11) is given by:

$$\min_{\boldsymbol{\pi} \in \mathbb{R}^{|\mathcal{I}|}} \qquad \sum_{i \in \mathcal{I}} \pi_i
\text{s.t.} \qquad -\lambda \le \sum_{i \in \mathcal{I}} y_i x_{ij} \pi_i \le \lambda \qquad j \in [p]
\qquad \sum_{i \in \mathcal{I}} y_i \pi_i = 0
0 \le \pi_i \le 1 \qquad i \in \mathcal{I}.$$
(12)

Let $(\beta^{\dagger}, \beta_0^{\dagger}) \in \mathbb{R}^{p+1}$ and $\pi^{\dagger} \in \mathbb{R}^{|\mathcal{I}|}$ denote optimal solutions of Problem (11) and (12) (respectively). Note that the reduced cost $\bar{\pi}_i$ of a dual variable π_i , $i \notin \mathcal{I}$ (cf. Section A.1) is

$$\bar{\pi}_i = 1 - y_i \left(\mathbf{x}_i^T \boldsymbol{\beta}^\dagger + \beta_0^\dagger \right).$$

We add to \mathcal{I} all indices corresponding to the dual variables having reduced cost higher than a threshold ϵ (specified a-priori). These indices correspond to violations of constraints (5a) in the primal Problem (5), and solve the new LP. The constraint generation algorithm for L1-SVM is summarized below:

Algorithm 3: Constraint generation for L1-SVM

Input: \mathbf{X} , \mathbf{y} , a regularization coefficient λ , a threshold $\epsilon > 0$, a set of constraints indexed by \mathcal{I} . Output: A near-optimal solution $\boldsymbol{\beta}^{\dagger}$ for the L1-SVM Problem (2).

- 1. Repeat Steps 2 to 3 until \mathcal{I} stabilizes.
- 2. Solve Problem $\mathcal{M}_{\ell_1}(\mathcal{I},[p])$ (i.e., Problem (11)).
- 3. Let $\mathcal{I}^{\epsilon} \subset \{1, \dots, n\} \setminus \mathcal{I}$ denote constraints with reduced cost higher than ϵ . Update $\mathcal{I} \leftarrow \mathcal{I} \cup \mathcal{I}^{\epsilon}$, and go to Step 2 (with LP warm-starting enabled).

Initialization: Similar to the case in column generation (cf Section 2.2.1), the constraint generation procedure benefits from a good initialization scheme. To this end, the first order methods described in Section 4.3 are found to be useful. A direct application of these first order methods suffer from increased computational cost when n becomes large due to the large cost associated with gradient computations. We thus use a (heuristic) sampling procedure that obtains approximate solutions (via a first order method) to the L1-SVM problem on different subsamples of the data and averages the estimators—leading to an estimate of the violated constraints \mathcal{I} (Section 4.4.2).

2.3.2 Column and constraint generation when both n and p are large

When both n and p are large, we will use a combination of column and constraint generation to solve the L1-SVM problem. For a given λ , let \mathcal{I} and \mathcal{J} denote subsets of columns and constraints (respectively). This leads to the following restricted version of the L1-SVM problem:

$$\begin{array}{ll}
\boxed{\mathcal{M}_{\ell_{1}}(\mathcal{I},\mathcal{J})} & \min_{\substack{\boldsymbol{\xi} \in \mathbb{R}^{|\mathcal{I}|}, \ \beta_{0} \in \mathbb{R} \\ \boldsymbol{\beta}^{+}, \ \boldsymbol{\beta}^{-} \in \mathbb{R}^{|\mathcal{J}|}}} & \sum_{i \in \mathcal{I}} \boldsymbol{\xi}_{i} + \lambda \sum_{j \in \mathcal{J}} \beta_{j}^{+} + \lambda \sum_{j \in \mathcal{J}} \beta_{j}^{-} \\ & \text{s.t.} & \boldsymbol{\xi}_{i} + \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{+} - \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{-} + y_{i} \beta_{0} \geq 1 \quad i \in \mathcal{I} \\
\boldsymbol{\xi} \geq 0, \ \boldsymbol{\beta}^{+} \geq 0, \ \boldsymbol{\beta}^{-} \geq 0.
\end{array} \tag{13}$$

Let $(\beta^*, \beta_0^*) \in \mathbb{R}^{|\mathcal{J}|+1}$, $\pi^* \in \mathbb{R}^{|\mathcal{I}|}$ be a pair of optimal primal and dual solutions for the above problem. Let $\bar{\beta}_j^+, \bar{\beta}_j^-$ denote the reduced costs for primal variables β_j^+, β_j^- and $\bar{\pi}_i$ denote the reduced cost for dual variable π_i . The reduced costs are given by:

$$\min\left\{\bar{\beta}_{j}^{+}, \bar{\beta}_{j}^{-}\right\} = \lambda - \left|\sum_{i \in \mathcal{I}} y_{i} x_{ij} \pi_{i}^{*}\right|; \quad \bar{\pi}_{i} = 1 - y_{i} \left(\sum_{j \in \mathcal{J}} x_{ij} \beta_{j}^{*} + \beta_{0}^{*}\right).$$
 (14)

We expand the sets \mathcal{I} and \mathcal{J} by using Steps 3 and 4 of Algorithm 4 (below). We then solve Problem (13) (with warm-starting enabled) and continue till \mathcal{I} and \mathcal{J} stabilize. Section 4.4.3 discusses the use of first order optimization methods to initialize \mathcal{I} and \mathcal{J} . Our hybrid column and constraint generation approach to solve the L1-SVM Problem (2) is summarized below.

Algorithm 4: Combined column and constraint generation for L1-SVM

Input: \mathbf{X} , \mathbf{y} , a regularization coefficient λ , a tolerance threshold $\epsilon > 0$, initial subsets \mathcal{I} and \mathcal{J} .

Output: A near-optimal solution $\boldsymbol{\beta}^*$ for the L1-SVM Problem (2).

- 1. Repeat Steps 2 to 4 until \mathcal{I} and \mathcal{J} stabilize.
- 2. Solve the model $\mathcal{M}_{\ell_1}(\mathcal{I},\mathcal{J})$ in Problem (13).
- 3. Let $\mathcal{I}^{\epsilon} \subset \{1, \dots, n\} \setminus \mathcal{I}$ denote constraints with reduced cost higher than ϵ . Update $\mathcal{I} \leftarrow \mathcal{I} \cup \mathcal{I}^{\epsilon}$.
- 4. Let $\mathcal{J}^{\epsilon} \subset \{1, \dots, p\} \setminus \mathcal{J}$ denote columns with reduced cost lower than $-\epsilon$. Update $\mathcal{J} \leftarrow \mathcal{J} \cup \mathcal{J}^{\epsilon}$; and go to Step 2.

2.4 Application to the Group-SVM problem

We now discuss how the framework presented above can be adapted to the Group-SVM Problem (3). We let $\mathcal{I}_g \subset [p]$ denote indices that belong to group g for $g \in [G]$.

Column generation: Below we present an LP formulation for Problem (3). We introduce the variables $\mathbf{v} = (v_q)_{q \in [G]}$ such that v_q refers to the L_{∞} -norm of the coefficients $\boldsymbol{\beta}_q$:

(Group-SVM)
$$\min_{\substack{\boldsymbol{\xi} \in \mathbb{R}^{n}, \beta_{0} \in \mathbb{R}, \\ \boldsymbol{\beta}^{+}, \boldsymbol{\beta}^{-} \in \mathbb{R}^{p}, \boldsymbol{v} \in \mathbb{R}^{G}}} \sum_{i=1}^{n} \xi_{i} + \lambda \sum_{g=1}^{G} v_{g}$$
s.t.
$$\xi_{i} + y_{i} \mathbf{x}_{i}^{T} \boldsymbol{\beta}^{+} - y_{i} \mathbf{x}_{i}^{T} \boldsymbol{\beta}^{-} + y_{i} \beta_{0} \geq 1 \quad i \in [n]$$

$$v_{g} - \beta_{j}^{+} - \beta_{j}^{-} \geq 0 \qquad j \in \mathcal{I}_{g}, \ g \in [G]$$

$$\boldsymbol{\xi} \geq 0, \ \boldsymbol{\beta}^{+} \geq 0, \ \boldsymbol{\beta}^{-} \geq 0, \ \boldsymbol{v} \geq 0.$$

$$(15)$$

A dual of Problem (15) is given by:

(Dual-Group-SVM)
$$\max_{\boldsymbol{\pi} \in \mathbb{R}^n} \qquad \sum_{i=1}^n \pi_i$$
s.t.
$$\sum_{j \in \mathcal{I}_g} \left| \sum_{i=1}^n y_i x_{ij} \pi_i \right| \le \lambda \qquad g \in [G]$$

$$\mathbf{y}^T \boldsymbol{\pi} = 0$$

$$0 \le \pi_i \le 1 \qquad i \in [n].$$
(16)

Following the description in Section A.1, we apply column generation on the groups. We bring into the model groups with lowest reduced cost. Here, the reduced cost of group g is given as:

$$\bar{\beta}_g = \lambda - \sum_{i \in \mathcal{I}_g} \left| \sum_{i=1}^n y_i x_{ij} \pi_i \right|. \tag{17}$$

Computing a regularization path: The regularization path algorithm presented in Section 2.2.2 can be adapted to the Group-SVM problem. First, note that:

$$\boldsymbol{\beta}^*(\lambda) = \mathbf{0}, \quad \forall \lambda \ge \lambda_{\max} = \max_{g \in [G]} \sum_{j \in \mathcal{I}_g} \sum_{i=1}^n |x_{ij}|. \tag{18}$$

For $\lambda = \lambda_{\text{max}}$, the reduced cost of variables corresponding to group g is given by the "group" analogue of (10):

$$\bar{\beta}_g = \lambda_{\max} - \sum_{j \in \mathcal{I}_g} \left| \frac{N_-}{N_+} \sum_{i \in \mathcal{I}_+} y_i x_{ij} + \sum_{i \in \mathcal{I}_-} y_i x_{ij} \right|. \tag{19}$$

As in Section 2.2.2, we can obtain a small set of groups maximizing the rhs of (19). We use these groups to initialize the LP solver to solve Problem (15) for the next small value of λ , using column generation—this results in computational savings when the number of active groups is small compared to G. We repeat this process for smaller values of λ using warm-start continuation.

Constraint generation and column generation: When n is large (but the number of groups is small) constraint generation can be used for the Group-SVM problem in a manner similar to that used for the L1-SVM problem. Similarly, column and constraint generation can be applied together to obtain computational savings when both n and the number of groups are large.

First order methods: First order methods (cf Section 4) can be used to obtain approximate solutions to the Group-SVM problem — they are found to be useful to obtain good initializations for the column and/or constraint generation methods.

3 Cutting plane algorithms for Slope-SVM

Here we discuss the Slope-regularized SVM estimator i.e., Problem (4). For a p-dimensional regularization parameter λ with coordinates sorted as: $\lambda_1 \geq \ldots \geq \lambda_p \geq 0$, we let

$$\|\boldsymbol{\beta}\|_{S} := \sum_{j=1}^{p} \lambda_{j} |\beta_{(j)}| \tag{20}$$

denote the Slope norm (for convenience, we drop the dependence on λ in the notation $\|\cdot\|_S$).

We note that the epigraph of $\|\beta\|_S$ i.e., $\{(\beta, \eta) \mid \|\beta\|_S \leq \eta\}$ can be expressed with exponentially many linear inequalities (Section 3.1) when using O(p)-many variables. We note that this epigraph admits an LP formulation using $O(p^2)$ many variables and $O(p^2)$ many constraints (cf Section A.2)—given the problem-sizes we seek to address, we do not pursue this route either. The large number of constraints makes column/constraint generation methodology for the Slope penalty significantly more challenging than the L1-SVM case. Section 3.1 discusses a constraint generation method that significantly reduces the number of constraints needed to model the epigraph. Section 3.2 discusses the use of column generation to exploit sparsity in β when p is large. Finally, Section 3.3 combines these two features to address the Slope SVM problem. We note that both column and constraint generation methods are needed for the Slope penalty, making it different from the L1-penalty, where column generation (alone) suffices. In what follows, we concentrate on the case where n is small but p is large — if n is also large, an additional layer of constraint generation might be needed to efficiently handle sparsity arising from the hinge-loss.

3.1 Constraint generation for Slope-SVM

Reformulation of Slope-SVM: Note that Problem (4) can be expressed as:

$$\begin{array}{ll}
\boxed{\mathcal{M}_{S}(\mathcal{C},[p])} & \underset{\boldsymbol{\xi} \in \mathbb{R}^{n}, \ \beta_{0}, \eta \in \mathbb{R}, \\ \boldsymbol{\beta}^{+}, \ \boldsymbol{\beta}^{-} \in \mathbb{R}^{p}} & \sum_{i=1}^{n} \xi_{i} + \eta \\
\text{s.t.} & \xi_{i} + y_{i} \mathbf{x}_{i}^{T} \boldsymbol{\beta}^{+} - y_{i} \mathbf{x}_{i}^{T} \boldsymbol{\beta}^{-} + y_{i} \beta_{0} \geq 1 \quad i \in [n] \\
& (\boldsymbol{\beta}^{+}, \ \boldsymbol{\beta}^{-}, \ \eta) \in \mathcal{C} \\
\boldsymbol{\xi} \geq 0, \ \boldsymbol{\beta}^{+} \geq 0, \ \boldsymbol{\beta}^{-} \geq 0
\end{array} \tag{21a}$$

where, $\beta = \beta^+ - \beta^-$; and β^+ , β^- denote the positive and negative parts of β (respectively). We express the Slope penalty in the epigraph form (21b) with \mathcal{C} defined as:

$$\mathcal{C} := \left\{ (\boldsymbol{\beta}^+, \ \boldsymbol{\beta}^-, \ \boldsymbol{\eta}) \ \middle| \ \boldsymbol{\eta} \geq \sum_{j=1}^p \lambda_j \beta_{(j)}^+ + \sum_{j=1}^p \lambda_j \beta_{(j)}^-, \ \boldsymbol{\beta}^+, \ \boldsymbol{\beta}^- \in \mathbb{R}^p \right\}$$

where, we use the notation $\beta_{(1)}^+ + \beta_{(1)}^- \ge ... \ge \beta_{(p)}^+ + \beta_{(p)}^-$ and remind ourselves that $|\beta_i| = \beta_i^+ + \beta_i^-$ for all *i*. Below we show that (21b) can be expressed via linear inequalities involving (β^+, β^-) .

Let S_p denote the set of all permutations of $\{1, \ldots, p\}$, with $|S_p| = p!$. For a permutation $\phi \in S_p$, we let $(\phi(1), \ldots, \phi(p))$ denote the corresponding rearrangement of $(1, \ldots, p)$. With this notation in place, note that the Slope norm can be expressed as:

$$\|\beta\|_{S} = \sum_{j=1}^{p} \lambda_{j} |\beta_{(j)}| = \max_{\phi \in \mathcal{S}_{p}} \sum_{j=1}^{p} \lambda_{j} |\beta_{\phi(j)}| = \max_{\psi \in \mathcal{S}_{p}} \sum_{j=1}^{p} \lambda_{\psi(j)} |\beta_{j}|.$$
 (22)

As a consequence, we have the following lemma:

Lemma 1. The Slope norm $\|\beta\|_S$ admits the following representation

$$\|\boldsymbol{\beta}\|_{S} = \max_{\boldsymbol{w} \in \mathcal{W}^{[p]}} \boldsymbol{w}^{T}(\boldsymbol{\beta}^{+} + \boldsymbol{\beta}^{-}) = \max_{\boldsymbol{w} \in \mathcal{W}^{[p]}_{o}} \boldsymbol{w}^{T}(\boldsymbol{\beta}^{+} + \boldsymbol{\beta}^{-})$$

where, $\mathcal{W}_0^{[p]} := \operatorname{Conv}\left(\mathcal{W}^{[p]}\right)$ is the convex hull of $\mathcal{W}^{[p]}$, where

$$\mathcal{W}^{[p]} := \left\{ \boldsymbol{w} \in \mathbb{R}^p \mid \exists \psi \in \mathcal{S}_p \quad s.t. \quad w_j = \lambda_{\psi(j)}, j \in [p] \right\}. \tag{23}$$

Proof. Note that a linear function maximized over a bounded polyhedron reaches its maximum at one of the extreme points of the polyhedron — this leads to:

$$\max_{\boldsymbol{w} \in \mathcal{W}_0^{[p]}} \boldsymbol{w}^T (\boldsymbol{\beta}^+ + \boldsymbol{\beta}^-) = \max_{\boldsymbol{w} \in \mathcal{W}^{[p]}} \boldsymbol{w}^T (\boldsymbol{\beta}^+ + \boldsymbol{\beta}^-).$$
 (24)

Using the definition of $\mathcal{W}^{[p]}$, we get that the rhs of (24) is $\max_{\psi \in \mathcal{S}_p} \sum_{j=1}^p \lambda_{\psi(j)} |\beta_j|$ which is in fact the Slope norm $\|\boldsymbol{\beta}\|_S$.

The following remark provides a description of $\mathcal{W}^{[p]}$ for certain choices of λ .

Remark 1. (a) If all the coefficients are equal i.e., $\lambda_1 = \ldots = \lambda_p$ and $\|\boldsymbol{\beta}\|_S = \lambda \|\boldsymbol{\beta}\|_1$, then $\mathcal{W}^{[p]}$ is a singleton. (b) If all the coefficients are distinct i.e., $\lambda_1 > \ldots > \lambda_p$, then each permutation $\psi \in \mathcal{S}_p$ is associated with a unique vector in $\mathcal{W}^{[p]}$ and $\mathcal{W}^{[p]}$ contains p! elements.

Using Lemma 1, we can derive an LP formulation of Problem (21) by modeling \mathcal{C} in (21b) as:

$$C = \left\{ \left(\boldsymbol{\beta}^+, \ \boldsymbol{\beta}^-, \ \eta \right) \ \middle| \ \boldsymbol{\beta}^+, \boldsymbol{\beta}^- \in \mathbb{R}^p, \ \eta \ge \max_{\boldsymbol{w} \in \mathcal{W}^{[p]}} \boldsymbol{w}^T (\boldsymbol{\beta}^+ + \boldsymbol{\beta}^-) \right\}, \tag{25}$$

where, $\mathcal{W}^{[p]}$ is defined in (23). The resulting LP formulation (21) has at most n constraints from (21a) and at most p! constraints associated with (21b) (by virtue of (25)). We note that many constraints in (25) are redundant: for example, the maximum is attained corresponding to the inverse of permutation ϕ (denoted by ϕ^{-1}), where $|\beta_{\phi(1)}| \geq \ldots \geq |\beta_{\phi(p)}|$. This motivates the use of constraint generation techniques.

Constraint generation: We proceed by replacing $\mathcal{W}^{[p]}$ with a smaller subset and solve the resulting LP. We subsequently refine this approximation if (21b) is violated. Formally, let us consider a collection of vectors/cuts $\mathbf{w}^{(1)}, \ldots, \mathbf{w}^{(t)} \in \mathcal{W}^{[p]}$ leading to a subset \mathcal{C}_t of \mathcal{C} :

$$\mathcal{C}_t := \left\{ (\boldsymbol{\beta}^+, \ \boldsymbol{\beta}^-, \ \eta) \ \middle| \ \boldsymbol{\beta}^+, \boldsymbol{\beta}^- \in \mathbb{R}^p, \ \eta \ge \sum_{j=1}^p w_j^{(\ell)} \beta_j^+ + \sum_{j=1}^p w_j^{(\ell)} \beta_j^-, \ \forall \ell \le t \right\} \subseteq \mathcal{C}.$$
 (26)

By replacing C in (21b) by C_t , we get an LP denoted by $\mathcal{M}_S(C_t, [p])$ which is a relaxation of $\mathcal{M}_S(C, [p])$. Let (β^*, η^*) be a solution of $\mathcal{M}_S(C_t, [p])$. If this is not an optimal solution (at a tolerance threshold $\epsilon > 0$), we add a cut to C_t if

$$\eta^* + \epsilon < \sum_{j=1}^p \lambda_j |\beta_{(j)}^*| = ||\beta^*||_S.$$

To this end, consider a permutation $\psi_{t+1} \in \mathcal{S}_p$ such that $|\beta_{\psi_{t+1}(1)}^*| \geq \ldots \geq |\beta_{\psi_{t+1}(p)}^*|$. If ψ_{t+1}^{-1} denotes the inverse of ψ_{t+1} , we obtain $\mathbf{w}^{(t+1)} \in \mathcal{W}^{[p]}$ such that:

$$w_j^{(t+1)} = \lambda_{\psi_{t+1}^{-1}(j)} \ \forall j \in [p]$$
 (27)

and solve the resulting LP. We continue adding cuts, till no further cuts need to be added — this leads to a (near)-optimal solution to $\mathcal{M}_S(\mathcal{C},[p])$. We note that the first cut $\mathbf{w}^{(1)}$ can be obtained by applying (27) on an estimator obtained from the first order optimization schemes (cf Section 4). Our algorithm is summarized below for convenience.

Algorithm 5: Constraint generation for Slope-SVM

Input: **X**, **y**, a vector of Slope coefficients $\{\lambda_j\}_{j\in[p]}$, a tolerance threshold $\epsilon > 0$, a cut $\mathbf{w}^{(1)} \in \mathcal{W}^{[p]}$. Output: A near-optimal solution $\boldsymbol{\beta}^*$ for the Slope-SVM Problem (4).

1. Repeat Steps 2 to 3 (for $t \ge 1$) till no further cuts need to be added.

- 2. Solve the model $\mathcal{M}_S(\mathcal{C}_t, [p])$ with \mathcal{C}_t as in (26). Let $(\boldsymbol{\beta}^*, \eta^*)$ be a solution.
- 3. Let $\psi_{t+1} \in \mathcal{S}_p$ be such that $|\beta^*_{\psi_{t+1}(1)}| \geq \ldots \geq |\beta^*_{\psi_{t+1}(p)}|$. If condition $\eta^* + \epsilon \geq \sum_{j=1}^p \lambda_j |\beta^*_{\psi_{t+1}(j)}|$ is not satisfied, we add a new cut $\mathbf{w}^{(t+1)} \in \mathcal{W}^{[p]}$ as per (27); update \mathcal{C}_{t+1} and go to Step 2.

3.2 Dual formulation and column generation for Slope-SVM

When the amount of regularization is high, the Slope penalty (with $\lambda_i > 0$ for all i) will lead to many zeros in an optimal solution to Problem (4) — computational savings are possible if we can leverage this sparsity when p is large. To this end, we use column generation along with the cutting plane algorithm described in Section 3.1. In particular, given a set of columns $\mathcal{J} = \{\mathcal{J}(1), \ldots, \mathcal{J}(|\mathcal{J}|)\} \subset \{1, \ldots, p\}$, we consider a restricted version of Problem (4) with $\beta_j = 0, j \notin \mathcal{J}$:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \beta_0 \in \mathbb{R}} \quad \sum_{i=1}^n \left(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \right)_+ + \|\boldsymbol{\beta}\|_S \quad \text{s.t.} \quad \beta_j = 0, j \notin \mathcal{J}.$$

The above can be expressed as an LP similar to Problem (21) but with a restricted set of columns

$$\begin{array}{ll}
\boxed{\mathcal{M}_{S}\left(\mathcal{C}^{\mathcal{J}},\mathcal{J}\right)} & \underset{\boldsymbol{\xi} \in \mathbb{R}^{n}, \ \beta_{0}, \ \eta \in \mathbb{R}, \\ \boldsymbol{\beta}_{\mathcal{J}}^{+}, \ \boldsymbol{\beta}_{\mathcal{J}}^{-} \in \mathbb{R}^{|\mathcal{J}|} & \sum_{i=1}^{n} \xi_{i} + \eta \\
& \text{s.t.} & \xi_{i} + \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{+} - \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{-} + y_{i} \beta_{0} \geq 1 & i \in [n] \\
& (\boldsymbol{\beta}_{\mathcal{J}}^{+}, \ \boldsymbol{\beta}_{\mathcal{J}}^{-}, \ \eta) \in \mathcal{C}^{\mathcal{J}} \\
& \boldsymbol{\xi} \geq 0, \ \boldsymbol{\beta}_{\mathcal{J}}^{+} \geq 0, \ \boldsymbol{\beta}_{\mathcal{J}}^{-} \geq 0
\end{array}$$

where, $\beta_{\mathcal{J}}$ is a sub-vector of β restricted to \mathcal{J} and $\mathcal{C}^{\mathcal{J}}$ is the adaption of (25) restricted to $\beta_{\mathcal{J}}$:

$$C^{\mathcal{J}} := \left\{ \left(\boldsymbol{\beta}_{\mathcal{J}}^{+}, \ \boldsymbol{\beta}_{\mathcal{J}}^{-}, \ \boldsymbol{\eta} \right) \ \middle| \ \boldsymbol{\beta}_{\mathcal{J}}^{+}, \boldsymbol{\beta}_{\mathcal{J}}^{-} \in \mathbb{R}^{|\mathcal{J}|}, \ \boldsymbol{\eta} \ge \max_{\boldsymbol{w}_{\mathcal{J}} \in \mathcal{W}^{\mathcal{J}}} \boldsymbol{w}_{\mathcal{J}}^{T} (\boldsymbol{\beta}_{\mathcal{J}}^{+} + \boldsymbol{\beta}_{\mathcal{J}}^{-}) \right\}$$
(29)

where, $\boldsymbol{w}_{\mathcal{J}} \in \mathbb{R}^{|\mathcal{J}|}$ and $\mathcal{W}^{\mathcal{J}}$ is defined as:

$$\mathcal{W}^{\mathcal{J}} := \left\{ oldsymbol{w}_{\mathcal{J}} \; \middle| \; \exists \psi \in \mathcal{S}_{|\mathcal{J}|} \; \; \mathrm{s.t.} \; \; w_{\mathcal{J}(j)} = \lambda_{\psi(j)}, \; \forall j \leq |\mathcal{J}| \right\}.$$

Since column generation is equivalent to constraint generation on the dual problem, to determine the set of columns to add to \mathcal{J} in Problem (28), we need the dual formulation of Slope-SVM.

Dual formulation for Slope-SVM: We first present a dual [37] of the Slope norm:

$$\max \left\{ \boldsymbol{\beta}^T \boldsymbol{z} \mid \boldsymbol{\beta} \in \mathbb{R}^p, \ \|\boldsymbol{\beta}\|_S \le 1 \right\} = \max_{k \le p} \left\{ \left(\sum_{j=1}^k \lambda_j \right)^{-1} \sum_{j=1}^k |z_{(j)}| \right\}.$$
 (30)

The identity (30) follows from the observation that the maximum will be attained at an extreme point of the polyhedron $\mathcal{P}_S = \{\beta \mid \|\beta\|_S \leq 1\} \subset \mathbb{R}^p$. We describe these extreme points. We fix $k \in [p]$, and a subset $A \subset \{1, \ldots, p\}$ of size k— the extreme points of \mathcal{P}_S having support A have their nonzero coefficients to be equal, with absolute value $\left(\sum_{j=1}^k \lambda_j\right)^{-1}$. Finally, (30) follows by taking a maximum over all $k \in [p]$.

A dual of Problem (28) is given by:

$$\max_{\boldsymbol{\pi} \in \mathbb{R}^{n}, \mathbf{q} \in \mathbb{R}^{p}} \qquad \sum_{i=1}^{n} \pi_{i}$$
s.t.
$$\max_{k=1,\dots,|\mathcal{J}|} \left\{ \left(\sum_{j=1}^{k} \lambda_{j} \right)^{-1} \sum_{j=1}^{k} |q_{(j)}| \right\} \leq 1$$

$$q_{j} = \sum_{i=1}^{n} y_{i} x_{ij} \pi_{i}, \quad j \in [p]$$

$$\mathbf{y}^{T} \boldsymbol{\pi} = 0$$

$$0 \leq \pi_{i} \leq 1, \quad i \in [n].$$

$$(31a)$$

We now discuss how additional columns can be appended to \mathcal{J} in Problem (28) to perform column generation. Let $\pi^* \in \mathbb{R}^n$ be an optimal solution of Problem (31). We compute the associated \mathbf{q}^* and sort its entries such that $|q_{(1)}^*| \geq \ldots \geq |q_{(|\mathcal{J}|)}^*|$. Constraint (31a) leads to:

$$\max_{k=1,\dots,|\mathcal{J}|} \left\{ \sum_{j=1}^{k} |q_{(j)}^*| - \sum_{j=1}^{k} \lambda_j \right\} \le 0.$$
 (32)

Now, for each column $j \notin \mathcal{J}$, we compute its corresponding $q_{(j)}^*$ and insert it into the sorted sequence $|q_{(1)}^*| \geq \ldots \geq |q_{(|\mathcal{J}|)}^*|$. This insertion costs at most $O(|\mathcal{J}|)$ flops: we update $\mathcal{J} \leftarrow \mathcal{J} \cup \{j\}$ and denote the sorted entries by: $|q^*_{(1)}| \geq \ldots \geq |q^*_{(|\mathcal{J}|+1)}|$. We add a column $j \notin \mathcal{J}$ to the current model if:

$$\max_{k=1,\dots,|\mathcal{J}|+1} \left\{ \sum_{j=1}^{k} |q^*_{\overline{(j)}}| - \sum_{j=1}^{k} \lambda_j \right\} > \epsilon, \tag{33}$$

and this costs $O(|\mathcal{J}|+1)$ flops. Therefore, the total cost of sorting the vector \mathbf{q}^* and scanning through all columns (not in the current model) for negative reduced costs, is of the order $O(|\mathcal{J}|\log|\mathcal{J}|+2(p-|\mathcal{J}|)|\mathcal{J}|)$. This approach can be computationally expensive. To this end, we propose an alternative method having a smaller cost with $O(|\mathcal{J}|)$ flops. Indeed, by combining

Equations (33) and (32), a column $j \notin \mathcal{J}$ will be added to the model if it satisfies:

$$|q_j| \ge \lambda_{|\mathcal{J}|+1} + \epsilon. \tag{34}$$

This shows that the cost of adding a new column for Slope-SVM is the same as that in L1-SVM. The column generation algorithm is summarized below.

Algorithm 6: Column generation for Slope-SVM

Input: \mathbf{X} , \mathbf{y} , a sequence of Slope coefficients $\{\lambda_j\}$, a threshold $\epsilon > 0$, an initial set of columns \mathcal{J} . Output: A near-optimal solution β^* for the Slope-SVM Problem (4).

- 1. Repeat Steps 2 to 3 until no column can be added.
- 2. Solve the model $\mathcal{M}_S(\mathcal{C}^{\mathcal{I}}, \mathcal{J})$ in Problem (26) with warm-start (if available).
- 3. Identify the columns $\mathcal{J}^{\epsilon} \subset \{1, \dots, p\} \setminus \mathcal{J}$ that need to be added by using criterion (34). Update $\mathcal{J} \leftarrow \mathcal{J} \cup \mathcal{J}^{\epsilon}$, and go to Step 2.

3.3 Pairing column and constraint generation for Slope-SVM

We discuss how to combine the column (Section 3.2) and constraint generation methods (Section 3.1) outlined above to solve the Slope SVM problem.

For a set of columns \mathcal{J} and constraints associated with $\boldsymbol{w}_{\mathcal{J}}^{(1)}, \dots, \boldsymbol{w}_{\mathcal{J}}^{(t)} \in \mathcal{W}^{\mathcal{J}}$, we consider the following problem

$$\begin{array}{c}
\left[\begin{array}{c} \mathcal{M}_{S}\left(\mathcal{C}_{t}^{\mathcal{J}},\mathcal{J}\right) \end{array}\right] & \underset{\boldsymbol{\xi} \in \mathbb{R}^{n}, \ \beta_{0} \in \mathbb{R}, \ \eta \in \mathbb{R}}{\min} \\ \boldsymbol{\beta}_{\mathcal{J}}^{+}, \ \boldsymbol{\beta}_{\mathcal{J}}^{-} \in \mathbb{R}^{|\mathcal{J}|} \end{array} \qquad \qquad \underset{i=1}{\overset{\sum}{\sum}} \xi_{i} + \eta \\
\text{s.t.} & \xi_{i} + \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{+} - \sum_{j \in \mathcal{J}} y_{i} x_{ij} \beta_{j}^{-} + y_{i} \beta_{0} \geq 1 \qquad i \in [n] \\
\left(\boldsymbol{\beta}_{\mathcal{J}}^{+}, \ \boldsymbol{\beta}_{\mathcal{J}}^{-}, \ \eta\right) \in \mathcal{C}_{t}^{\mathcal{J}} \\
\boldsymbol{\xi} \geq 0, \ \boldsymbol{\beta}_{\mathcal{J}}^{+} \geq 0, \ \boldsymbol{\beta}_{\mathcal{J}}^{-} \geq 0,
\end{array}$$

where, $C_t^{\mathcal{J}}$ (and $C^{\mathcal{J}}$) are restrictions of C_t (and C, respectively) to the columns \mathcal{J} . Formally,

$$\mathcal{C}_t^{\mathcal{J}} := \left\{ \left(oldsymbol{eta}_{\mathcal{J}}^+, \; oldsymbol{eta}_{\mathcal{J}}^-, \; \eta
ight) \; \middle| \; oldsymbol{eta}_{\mathcal{J}}^+, \; oldsymbol{eta}_{\mathcal{J}}^- \in \mathbb{R}^{|\mathcal{J}|}, \quad \eta \geq (oldsymbol{w}_{\mathcal{J}}^{(\ell)})^T (oldsymbol{eta}_{\mathcal{J}}^+ + oldsymbol{eta}_{\mathcal{J}}^-), \; orall \ell \leq t
ight\} \subset \mathcal{C}^{\mathcal{J}}.$$

We use the method in Section 3.1 to refine $C_t^{\mathcal{J}}$ and the method of Section 3.2 to add a set of columns to \mathcal{J} . We use criterion (34) to select the columns to add. Let \mathcal{J}^{ϵ} denote these additional columns with coordinates $\mathcal{J}^{\epsilon}(k)$ for $k = 1, \ldots, |\mathcal{J}^{\epsilon}|$ — we will also assume that the elements have

been sorted by increasing reduced costs. For notational purposes, we will need to map⁴ the existing cuts of $W^{\mathcal{J}}$ onto $W^{\mathcal{J}\cup\mathcal{J}^{\epsilon}}$. To this end, we make the following definition:

$$w_m^{(\ell)} = \lambda_{|\mathcal{J}|+k}, \ \forall m \in \mathcal{J}^{\epsilon}, \ \forall \ell \le t.$$
 (36)

We summarize our algorithm below.

Algorithm 7: Column-and-constraint generation for Slope-SVM

Input: \mathbf{X} , \mathbf{y} , a sequence of Slope coefficient $\{\lambda_j\}$, a threshold $\epsilon > 0$. Initialization of $\boldsymbol{\beta}^*$ and $\boldsymbol{\mathcal{J}}$ (e.g., using the first order method in Section 4.3). Define $\boldsymbol{w}_{\mathcal{J}}^{(1)}$ as per (27).

Output: A near-optimal solution β^* for the Slope-SVM Problem (4).

- 1. Repeat Steps 2 to 4 until no cut can be added and $\mathcal J$ stabilizes.
- 2. Solve the model $\mathcal{M}_S\left(\mathcal{C}_t^{\mathcal{I}},\mathcal{J}\right)$ in Problem (35) (with warm-starting enabled).
- 3. If $\eta < \sum_{j=1}^{|\mathcal{J}|} \lambda_j |\beta_{(j)}^*| \epsilon$, add a new cut $\boldsymbol{w}_{\mathcal{J}}^{(t+1)} \in \mathcal{W}^{\mathcal{J}}$ as in Equation (27) and define $\mathcal{C}_{t+1}^{\mathcal{J}}$.
- 4. Identify columns $\mathcal{J}^{\epsilon} \subset \{1, \dots, p\} \setminus \mathcal{J}$ that need to be added (based on criterion (34)). Map the cuts $\boldsymbol{w}_{\mathcal{J}}^{(1)}, \dots \boldsymbol{w}_{\mathcal{J}}^{(t+1)}$ to $\mathcal{W}^{\mathcal{J} \cup \mathcal{J}^{\epsilon}}$ via (36). Update $\mathcal{J} \leftarrow \mathcal{J} \cup \mathcal{J}^{\epsilon}$ and go to Step 2.

4 First order methods

The computational performance of the cutting plane methods described above (for Problems (2), (3), (4)), are found to benefit from a good initialization (e.g., a good estimate of columns for column generation) especially, when compared to initializing with a random set of columns.

For initialization purposes, we use a low-accuracy solution⁵ obtained via first order methods [23]. Since all the problems ((2), (3), (4)) are nonsmooth, we use Nesterov's smoothing technique [24] to smooth the nonsmooth hinge-loss function and use proximal gradient descent on the composite version [25] of the problem⁶. These solutions serve as reasonable (initial) estimates for the sets of columns (respectively constraints) necessary for the column (respectively constraint) generation methods. When the number of samples and/or features become larger, a direct application of the first order methods becomes expensive and we use additional heuristics (e.g., correlation screening, sub-sampling) for scalability (Section 4.4).

⁴In other words, the existing vectors $\mathbf{w}_{\mathcal{J}}^{(\ell)}$ are in $\mathbb{R}^{|\mathcal{J}|}$ and we need to extend them to $\mathbb{R}^{|\mathcal{J}|+|\mathcal{J}^{\epsilon}|}$. Therefore, we need to define the coordinates corresponding to the new indices \mathcal{J}^{ϵ} .

⁵Obtaining high accuracy solutions via first order methods can become prohibitively expensive especially, when compared to the cutting plane algorithms presented here.

⁶For the Group-SVM problem, we use proximal block coordinate methods instead of proximal gradient methods.

4.1 Solving the composite form with Nesterov's smoothing

Note that for a scalar u, we have $\max\{0, u\} = \frac{1}{2}(u + |u|) = \max_{|w| \le 1} \frac{1}{2}(u + wu)$ and this maximum is achieved when w = sign(x). Hence, the hinge-loss can be expressed as:

$$\sum_{i=1}^{n} \left(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) \right)_+ = \max_{\|\boldsymbol{w}\|_{\infty} \le 1} \sum_{i=1}^{n} \frac{1}{2} \left[1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) + w_i(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0)) \right], \tag{37}$$

and its smoothed version [24] is given by:

$$F^{\tau}(\boldsymbol{\beta}, \beta_0) = \max_{\|\boldsymbol{w}\|_{\infty} \le 1} \left\{ \sum_{i=1}^{n} \frac{1}{2} \left[1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0) + w_i(1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0)) \right] - \frac{\tau}{2} \|\boldsymbol{w}\|_2^2 \right\}$$

where, $F^{\tau}(\boldsymbol{\beta}, \beta_0)$ is differentiable w.r.t $(\boldsymbol{\beta}, \beta_0)$. $F^{\tau}(\boldsymbol{\beta}, \beta_0)$ is a pointwise $O(\tau)$ -approximation to the hinge-loss (37). Using the notation: $\mathbf{z}, \mathbf{w}^{\tau} \in \mathbb{R}^n$: $z_i = 1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0)$, $\forall i$ and $w_i^{\tau} = \min(1, \frac{1}{2\tau}|z_i|) \operatorname{sign}(z_i)$, $\forall i$; the gradient of F^{τ} is given by:

$$\nabla F^{\tau}(\boldsymbol{\beta}, \ \beta_0) = -\frac{1}{2} \sum_{i=1}^{n} (1 + w_i^{\tau}) y_i \tilde{\mathbf{x}}_i \in \mathbb{R}^{p+1}, \tag{38}$$

where, $\tilde{\mathbf{x}}_i = (\mathbf{x}_i, 1) \in \mathbb{R}^{p+1}$. Note $\nabla F^{\tau}(\boldsymbol{\beta}, \beta_0)$ is Lipschitz continuous (cf. Theorem 1 in [24]):

$$\|\nabla F^{\tau}(\beta, \beta_0) - \nabla F^{\tau}(\beta', \beta'_0)\|_2 \le C^{\tau} \|(\beta, \beta_0) - (\beta', \beta'_0)\|_2$$

with parameter $C^{\tau} = \sigma_{\max}(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})/4\tau$ where, $\tilde{\mathbf{X}}_{n \times (p+1)}$ is a matrix with *i*th row $\tilde{\mathbf{x}}_i$. We use a proximal gradient method [3] to the following composite form of the smoothed-hinge-loss SVM problem with regularizer $\Omega(\boldsymbol{\beta})$

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \beta_0 \in \mathbb{R}} F^{\tau}(\boldsymbol{\beta}, \beta_0) + \Omega(\boldsymbol{\beta}), \tag{39}$$

where, $\Omega(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|_1$ for L1-SVM, $\Omega(\boldsymbol{\beta}) = \lambda \sum_{g=1}^G \|\boldsymbol{\beta}_g\|_{\infty}$ for Group-SVM and $\Omega(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_S$ for Slope-SVM. For these choices, the proximal/thresholding operators can be computed easily.

4.2 Thresholding operators

For notational convenience we set $\gamma = (\beta, \beta^0) \in \mathbb{R}^{p+1}$. Following [23, 3] for $L \geq C^{\tau}$, we have that $\gamma \mapsto Q_L(\gamma; \alpha)$ is an upper bound to $\gamma \mapsto F^{\tau}(\gamma)$, i.e, for all $\alpha, \gamma \in \mathbb{R}^{p+1}$:

$$F^{\tau}(\gamma) \le Q_L(\gamma; \alpha) := F^{\tau}(\alpha) + \nabla F^{\tau}(\alpha)^T (\gamma - \alpha) + \frac{L}{2} \|\gamma - \alpha\|_2^2.$$
 (40)

The proximal gradient method requires solving the following problem:

$$\hat{\gamma} = \underset{\gamma}{\operatorname{arg\,min}} \left\{ Q_L(\gamma; \boldsymbol{\alpha}) + \Omega(\boldsymbol{\beta}) \right\} = \underset{\gamma}{\operatorname{arg\,min}} \frac{1}{2} \left\| \gamma - \left(\boldsymbol{\alpha} - \frac{1}{L} \nabla F^{\tau}(\boldsymbol{\alpha}) \right) \right\|_2^2 + \frac{1}{L} \Omega(\gamma). \tag{41}$$

We denote: $\hat{\gamma} = (\hat{\beta}, \hat{\beta}^0)$. Note that $\hat{\beta}_0$ is simple to compute and $\hat{\beta}$ can be computed via the following thresholding operator (with $\mu > 0$):

$$S_{\mu\Omega}(\boldsymbol{\eta}) := \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{arg\,min}} \frac{1}{2} \|\boldsymbol{\beta} - \boldsymbol{\eta}\|_2^2 + \mu\Omega(\boldsymbol{\beta}). \tag{42}$$

Computation of the thresholding operator is discussed below for specific choices of Ω .

Thresholding operator when $\Omega(\beta) = ||\beta||_1$: In this case, $S_{\mu\Omega}(\eta)$ is available via componentwise soft-thresholding where, the scalar soft-thresholding operator is given by:

$$\underset{u \in \mathbb{R}}{\operatorname{arg \, min}} \ \frac{1}{2} (u - c)^2 + \mu |u| = \operatorname{sign}(c)(|c| - \mu)_+.$$

Thresholding operator when $\Omega(\beta) = \sum_{g \in [G]} \|\beta_g\|_{\infty}$: We first consider the projection operator that projects onto an L1-ball with radius μ

$$\tilde{\mathcal{S}}_{\frac{1}{\mu}\|\cdot\|_1}(\boldsymbol{\eta}) := \underset{\boldsymbol{\beta}}{\operatorname{arg\,min}} \ \frac{1}{2} \|\boldsymbol{\beta} - \boldsymbol{\eta}\|_2^2 \quad \text{s.t.} \quad \frac{1}{\mu} \|\boldsymbol{\beta}\|_1 \le 1.$$
(43)

From standard results pertaining to the Moreau decomposition [22] (see also [1]) we have:

$$S_{\mu\|\cdot\|_{\infty}}(\boldsymbol{\eta}) + \tilde{S}_{\frac{1}{\mu}\|\cdot\|_{1}}(\boldsymbol{\eta}) = \boldsymbol{\eta}$$

$$\tag{44}$$

for any η . Note that $\tilde{\mathcal{S}}_{\frac{1}{\mu}\|\cdot\|_1}(\eta)$ can be computed via a simple sorting operation [32, 33], leading to a solution for $\mathcal{S}_{\mu\|\cdot\|_{\infty}}(\eta)$. This observation can be used to solve Problem (42) with the Group-SVM regularizer by noticing that the problem separates across the G groups.

Thresholding operator when $\Omega(\beta) = \sum_{i \in [p]} \lambda_i |\beta_{(i)}|$: For the Slope regularizer, Problem (42) reduces to the following optimization problem:

$$\min_{\beta} \ \frac{1}{2} \|\beta - \eta\|_2^2 + \mu \sum_i \lambda_i |\beta_{(i)}|. \tag{45}$$

As noted by [7], at an optimal solution to Problem (45), the signs of β_j and η_j are the same. In addition, since λ_i 's are decreasing, a solution to Problem (45) can be found by solving the following

close relative to the isotonic regression problem [28]

$$\min_{\mathbf{u}} \quad \frac{1}{2} \|\mathbf{u} - \tilde{\boldsymbol{\eta}}\|_{2}^{2} + \sum_{j=1}^{p} \mu \lambda_{j} u_{j} \quad \text{s.t.} \quad u_{1} \ge \dots \ge u_{p} \ge 0$$
 (46)

where, $\tilde{\boldsymbol{\eta}}$ is a decreasing re-arrangement of the absolute values of $\boldsymbol{\eta}$, with $\tilde{\eta}_i \geq \tilde{\eta}_{i+1}$ for all i. If $\hat{\mathbf{u}}$ is a solution to Problem (46)—then its ith coordinate \hat{u}_i corresponds to $|\hat{\beta}_{(i)}|$ where, $\hat{\boldsymbol{\beta}}$ is an optimal solution of Problem (45).

4.3 Deterministic first order algorithms

Accelerated gradient descent: Let us denote the mapping (41) $\alpha \mapsto \hat{\gamma}$ by the operator: $\hat{\gamma} := \Theta(\alpha)$. The basic version of the proximal gradient descent algorithm performs the updates: $\alpha_{T+1} = \Theta(\alpha_T)$ (for $T \ge 1$) after starting with $\alpha_1 = (\beta_1, \beta_1^0)$. The accelerated gradient descent algorithm [3], which enjoys a faster convergence rate performs updates with a minor modification. It starts with $\alpha_1 = \tilde{\alpha}_0$, $q_1 = 1$ and then performs the updates: $\tilde{\alpha}_{T+1} = \Theta(\alpha_T)$ where, $\alpha_{T+1} = \tilde{\alpha}_T + \frac{q_T-1}{q_{T+1}}(\tilde{\alpha}_T - \tilde{\alpha}_{T-1})$ and $q_{T+1} = (1 + \sqrt{1 + 4q_T^2})/2$. This algorithm requires $O(1/\epsilon)$ iterations to reach an ϵ -optimal solution for the original problem (with the hinge-loss). We perform these updates till some tolerance criterion is satisfied, for example, $\|\alpha_{T+1} - \alpha_T\| \le \eta$ for some tolerance level $\eta > 0$. In most of our examples (cf Section 5), we set a generous tolerance of $\eta = 10^{-3}$ or run the algorithm with a limit on the total number of iterations (usually a couple of hundred)⁷.

Block Coordinate Descent (CD) for the Group-SVM problem: We describe a cyclical proximal block coordinate (CD) descent algorithm [35] for the smooth hinge-loss function with the group regularizer. This method exhibits superior numerical performance when compared to a full gradient descent algorithm. We note that [27] explore block CD like algorithms for a different class of group Lasso type problems⁸ and our approaches differ.

We perform a proximal gradient step on the gth group of coefficients (with all other blocks and β_0 held fixed) via:

$$\boldsymbol{\beta}_g^{t+1} \in \underset{\boldsymbol{\beta}_g}{\operatorname{arg\,min}} \frac{1}{2} \left\| \boldsymbol{\beta}_g - \boldsymbol{\beta}_g^t - \frac{1}{C_g^{\tau}} \left\{ \nabla F^{\tau}(\boldsymbol{\beta}_1^{t+1}, \dots \boldsymbol{\beta}_{g-1}^{t+1}, \boldsymbol{\beta}_g^t, \dots \boldsymbol{\beta}_G^t, \boldsymbol{\beta}_0^t) \right\}_{\mathcal{I}_g} \right\|_2^2 + \frac{\lambda}{C_g^{\tau}} \|\boldsymbol{\beta}_g\|_{\infty}, \quad (47)$$

where $\{\nabla F^{\tau}(\cdot)\}_{\mathcal{I}_g}$ denotes the gradient restricted to the coordinates \mathcal{I}_g and C_g^{τ} is its associated Lipschitz constant: $C_g^{\tau} = \sigma_{\max}(\mathbf{X}_{\mathcal{I}_g}^T\mathbf{X}_{\mathcal{I}_g})/4\tau$. We cyclically update the coefficients across each group $g \in [G]$ and then update β^0 . This continues till some convergence criterion is met.

⁷This choice is user-dependent, with an obvious tradeoff between computation time and the quality of solution.

⁸The authors in [27] consider a different class of problems than those studied here; and use exact minimization for every block (they use a squared error loss function).

Important computational savings are possible for this CD algorithm by a careful accounting of flops. As one moves from one group to the next, the whole gradient can be updated easily. Note that the gradient $\nabla F^{\tau}(\beta, \beta_0)$ restricted to block g is given by:

$$\left\{\nabla F^{\tau}(\boldsymbol{\beta}, \beta_0)\right\}_{\mathcal{I}_g} = -\frac{1}{2}\mathbf{X}_{\mathcal{I}_g}^T \left\{\mathbf{y} \circ (1 + \mathbf{w}^{\tau})\right\},$$

where 'o' denotes element-wise multiplication. If \mathbf{w}^{τ} is known, the above computation requires $n|\mathcal{I}_g|$ flops. Recall that \mathbf{w}^{τ} depends upon $\boldsymbol{\beta}$ via: $w_i^{\tau} = \min\left(1, \frac{1}{2\tau}|z_i|\right) \operatorname{sign}(z_i)$ where $z_i = 1 - y_i(\mathbf{x}_i^T \boldsymbol{\beta} + \beta_0)$, $\forall i$. If $\boldsymbol{\beta}$ changes from $\boldsymbol{\beta}^{\text{old}}$ to $\boldsymbol{\beta}^{\text{new}}$, then \mathbf{w}^{τ} changes via an update in $\mathbf{X}\boldsymbol{\beta}$ — this change can be efficiently computed by noting that: $\mathbf{X}\boldsymbol{\beta}^{\text{new}} = \sum_{g \in [G]} \mathbf{X}_{\mathcal{I}_g} \boldsymbol{\beta}_g^{\text{new}} = \mathbf{X}\boldsymbol{\beta}^{\text{old}} + \mathbf{X}_{\mathcal{I}_g} \Delta \boldsymbol{\beta}_g$ where, $\Delta \boldsymbol{\beta}_g = \boldsymbol{\beta}_g^{\text{new}} - \boldsymbol{\beta}_g^{\text{old}}$ is a change that is only restricted to block g. Hence updating \mathbf{w}^{τ} also requires $n|\mathcal{I}_g|$ operations. The above suggests that one sweep of block CD across all the coordinates has a cost similar to that of computing a full gradient. In addition, the following characteristics [15] lead to further computational savings: (a) Active set strategy: We anticipate many of the groups to be zero in an optimal solution. If a group which was at zero stays at zero, no additional computations are needed. (b) Warm start continuation: The CD algorithm is well suited for computing a path of solutions in λ via warm-starts. Even if we desire a solution for a single value of λ , a continuation strategy can be used to speed up the overall algorithm (See Section 5).

4.4 Scalability heuristics for large problem instances

When n and/or p becomes large, the first order algorithms become expensive. Since our goal is to get a *ball park* estimate of the initial columns and/or constraints for the cutting plane algorithms; for scalability purposes, we use principled heuristics (motivated by statistical principles) as a wrapper around the first order methods.

4.4.1 Correlation screening when p is large and n is small

When $p \gg n$, we use correlation screening [30] to restrict the number of features (or groups in the case of Group-SVM). We apply the first order methods on this reduced set of features. Usually, for L1-SVM and Slope-SVM, we select the top 10n columns with highest absolute inner product (note that the features are standardized to have unit L2-norm) with the output. For the Group-SVM problem: for each group, we compute the inner products between every feature within this group and the response, and take their L1-norm. We then sort these numbers and take the top n groups.

4.4.2 A subsampling heuristic when n is large and p is small

The methods described in Section 4.3 become expensive due to gradient computations when n becomes large. When n is large but p is small, we use a subsampling method inspired by [21]. To get an approximate solution to Problem (2) we apply the algorithm in Section 4.3 on a subsample $(y_i, \mathbf{x}_i), i \in \mathcal{A}$ with sample-indices $\mathcal{A} \subset [n]$. We (approximately) solve Problem (2) with $\lambda \leftarrow \frac{|\mathcal{A}|}{n}\lambda$ (to adjust the dependence of λ on the sample size) by using the algorithms in Section 4.3. Let the solution obtained be given by $\hat{\beta}(\mathcal{A})$. We obtain $\hat{\beta}(\mathcal{A}_j)$ for different subsamples $\mathcal{A}_j, j \in [Q]$ and average the estimators $\bar{\beta}_Q = \frac{1}{Q} \sum_{j \in [Q]} \hat{\beta}(\mathcal{A}_j)$. We maintain a counter for Q, and stop as soon as the average stabilizes $\bar{\beta}_Q$ i.e., $\|\bar{\beta}_Q - \bar{\beta}_{Q-1}\| \leq \mu_{\text{Tol}}$ for some tolerance threshold μ_{Tol} . The estimate $\bar{\beta}_Q$ is used to obtain the violated constraints for the SVM problem and serves to initialize the constraint generation method.

4.4.3 A subsampling heuristic when both n and p are large

Here we basically apply a combination of the ideas described above for large p (small n) and large n (small p). More specifically, we choose a subsample \mathcal{A}_j and for this subsample, we use correlation screening to reduce the number of features and obtain an estimator $\hat{\beta}(\mathcal{A}_j)$. We then average these estimators (across \mathcal{A}_j s) to obtain $\bar{\beta}_Q$. If the support of $\bar{\beta}_Q$ is too large, we sort the absolute values of the coefficients and retain the top few hundred coefficients (in absolute value) to initialize the column generation method. The estimator $\bar{\beta}_Q$ is used to identify the samples for which the hinge-loss is nonzero — these indices are used to initialize the constraint generation method.

Details regarding the subsample sizes, number of continuation steps, the number of thresholded coefficients, etc, appear in Section 5.1.4.

5 Numerical experiments

We demonstrate the performance of our different methods on synthetic and real datasets for varying n, p values. We use the LP solver of Gurobi version 6.5.2 with Python interface. All computations are performed in Python. All computations (unless otherwise specified) are carried out on a Mac-Book with processor 2.7 GHz 12-Core Intel Xeon E5, 64GB of RAM.

Sections 5.1, 5.2 and 5.3 present computational results for the L1-SVM, Group-SVM and Slope-SVM problems (respectively).

⁹We note that the estimates $\hat{\boldsymbol{\beta}}(\mathcal{A}_j)$ can all be computed in parallel.

¹⁰We note that when n is large (but p is not too large), basic principles of statistical inference [31] suggest that the estimators $\hat{\beta}(A_j)$ will be reasonable approximations to a minimizer of Problem (2) — we average the estimators for variance reduction.

5.1 Computational results for L1-SVM

We present herein our computational experience with regard to the L1-SVM problem.

5.1.1 Experiments on synthetic datasets for p large, n small

Data Generation: We consider n samples from a multivariate Gaussian distribution with covariance matrix $\Sigma = ((\sigma_{ij}))$ with $\sigma_{ij} = \rho$ if $i \neq j$ and $\sigma_{ij} = 1$ otherwise. Half of the samples are from the +1 class and have mean $\mu_+ = (\mathbf{1}_{k_0}, \mathbf{0}_{p-k_0})$. The other half are from the -1 class and have mean $\mu_- = -\mu_+$. We standardize the columns of \mathbf{X} to have unit L2-norm.

Computing a regularization path: We first consider the task of solving the L1-SVM problem for a sequence of λ values, leading to a regularization path. Our goal is to demonstrate the effectiveness of using column generation over approaches that do not use column generation. We fix $n = 100, k_0 = 10, \rho = 0.1$ and consider different values of p (with $p \gg n$) — a regime where column generation is likely to be useful. For each training set, we use a geometric sequence of 20 decreasing values of λ with common ratio 0.7. We compare across the following methods:

- "CLG": This is Algorithm 2 for computing a regularization path (cf Section 2.2.2) with column generation ¹¹. We take $j_0 = 10$ and different tolerance values $\epsilon \in \{0.5, 0.1, 0.01\}$.
- "LP wo warm-start": this solves Problem (2) using Gurobi's LP solver (no column generation is used) for different λ-values with no warm-starts across the λ-path.
- "LP warm-start": this solves Problem (2) using Gurobi's LP solver (no column generation is used), across different λ-values using LP warm-starts across the λ-path.

The default settings of Gurobi's LP solver is used in all examples. The results are presented in Table 1. For an algorithm 'Alg', a regularization parameter λ , and a replication rep $\in \{1, \ldots, R\}$, we let $f_{\lambda}^{\text{Alg}}(\text{rep})$ be the objective value of the method (for the unconstrained problem (2)) and $f_{\lambda}^{*}(\text{rep})$ be the lowest objective value among all methods. We define the averaged relative accuracy (ARA) of method 'Alg' (for a fixed value of λ) over all the R simulations as:

$$ARA(Alg, \lambda) = \frac{1}{R} \sum_{\text{rep}=1}^{R} \frac{f_{\lambda}^{Alg}(\text{rep}) - f_{\lambda}^{*}(\text{rep})}{f_{\lambda}^{*}(\text{rep})}.$$

As different algorithms use different metrics for convergence, to make the algorithms comparable, we use the measure 'ARA'. Table 1 shows the mean and standard deviation (in parenthesis) of the R = 10 repetitions. We consider the *total* time to compute the entire path and also the averaged

 $^{^{-11}}$ First order methods are *not* used as initialization. Since we compute a path of solutions for Problem (2), the LPs are warm-started via continuation across λ -values.

ARA for all values of λ . The results in Table 1 show that the proposed column generation method (with warm-starting enabled across the regularization path) leads to the best overall performance in terms of obtaining a good solution with the smallest run times.

Training times for computing a regularization path for L1-SVM with $p \gg n$

	p = 1,000		p = 10,000		p = 100,000	
Method	Time (s)	ARA (%)	Time (s)	ARA (%)	Time (s)	ARA (%)
		LP without	column genera	ation		
LP wo warm-start	77.2 (1.3)	0.0	777.4(2.9)	0.0	> 7200	0.0
LP warm-start	4.4 (0.1)	0.0	44.4 (0.6)	0.0	471.4 (10)	0.0
LP with column generation						
CLG, ϵ =0.5	0.15 (0.02)	4.8 (0.9)	0.36(0.03)	2.8 (0.8)	2.89 (0.24)	3.2 (1.1)
CLG, ϵ =0.1	0.25 (0.02)	0.1(0.1)	0.53 (0.05)	0.1(0.1)	3.98 (0.84)	0.1(0.0)
CLG, ϵ =0.01	0.28 (0.02)	0.0 (0.0)	0.61 (0.04)	0.0 (0.0)	4.74 (0.92)	0.0 (0.0)

Table 1: Times (in secs) for computing the L1-SVM problem with 20 values of the regularization parameter: we compare Gurobi's LP solver with and without warm-start; versus our proposed column generation (CLG) method with three tolerance levels (denoted by ϵ). The basic LP solver (without column generation) benefits from using warm-starting across λ values — without warm-starts the algorithm takes more than 2 hrs to converge for $p = 10^5$ — here, CLG leads to more than a 2,500-fold improvement in run time. Overall CLG leads to significant improvements. The LP on the full problem reaches the best objective values (smallest ARA), but CLG is also quite accurate, as evidence from the ARA values (reported in %).

Results for a fixed λ : We study the performance of different methods for a fixed value of the regularization parameter, set to $\lambda = 0.01 \lambda_{\text{max}}^{12}$. We compare across the following 5 methods:

- (a) "RP CLG": We compute a solution to Problem (2) at the desired value of λ , using a regularization path (RP) (aka continuation) approach. We compute solutions on a grid of 7 regularization parameter values in the range $\left[\frac{1}{2}\lambda_{\max},\lambda\right]$ using the column generation algorithm (cf Section 2.2.1) with a tolerance threshold $\epsilon=0.01$ for every value of λ .
- (b) "F0+CLG": This is the column generation method initialized with a first order (FO) method (cf Section 4.3) with smoothing parameter $\tau = 0.2$. We use a termination criterion of $\eta = 10^{-3}$ or a maximum number of $T_{\text{max}} = 200$ iterations for the FO method. We use correlation screening to retain the top 10n features before applying the FO method. Column generation uses a tolerance level of $\epsilon = 0.01$. The time displayed includes the time taken to run the FO method. For reference, we report the time taken to run column generation excluding the time of the first order method: "CLG wo FO".
- (c) "Cor. screening": This initializes the column generation method by using correlation screening to retain the top 50 features.

¹²This choice leads to a model that is moderately sparse, with $\lambda = \lambda_{\text{max}}$ corresponding to a null model. The number of nonzeros in an optimal solution β^* , are 60, 66, 69 corresponding to the values p = 20K, 50K, 100K respectively (recall that n = 100 in this example).

- (d) "Random init.": This initializes the column generation method with a random (chosen uniformly) subset of 50 features.
- (e) "LP solver": This is Gurobi's LP solver to solve the full LP model (without column generation).

Among the above: The comparative timings among (b), (c) and (d) show the importance of having a good initialization and in particular, the effectiveness of using a first order method to initialize the column generation method. Method (a) computes a regularization path (via column generation) to arrive at the desired value of λ — it does not use any first order method like (b) — thus any timing difference between (a) and (b) is due to the role played by the first order methods for warm-starting. Figure 1 shows the results for a synthetic dataset with n = 100, $k_0 = 10$, $\rho = 0.1$ and R = 10 replications (error bars show the standard deviation of values across the 10 replications).

It appears that as long as $p \leq 3,000$ all methods perform similarly. The run time for the basic LP solver increases as p increases. Column generation with random initialization as well as correlation screening improves over this basic LP solver — correlation screening shows marginal improvement over random initialization. Column generation is found to benefit the most when initialized with the first order method. The runtimes of the first order method is also quite small (compare CLG wo FO and FO+CLG) and further reduction in computation time is possible with a more efficient implementation of the first order method.

5.1.2 Real datasets

Results for a fixed λ : We consider four real gene expression microarray datasets¹³ each having a small n and a large p. We compare our best method "F0+CLG" (method (b)) (as evidenced from the experiments with synthetic datasets) with the "LP solver" (method (e)). All methods are run with specifications explained in Section 5.1.1: for "F0+CLG" once a solution from the first order method was obtained, we sorted the coefficients in absolute value and retained the top 100 coefficients to initialize the column generation method. We use a tolerance level $\epsilon = 10^{-2}$.

For each dataset we merged the training and test sets and perform an L2-norm standardization of every column of the feature matrix. Table 2 presents the averaged training times over 10 repetitions for $\lambda = 0.01 \lambda_{\rm max}$ — results show that our column generation method can lead to a 70-fold speedup over LP solver — the overall run times are small as n is small.

¹³The leukemia, lung cancer and ovarian datasets can be found at: http://cilab.ujn.edu.cn/Dataset.htm. The Radsens dataset is available at: https://statweb.stanford.edu/~tibs/ElemStatLearn/datasets/.

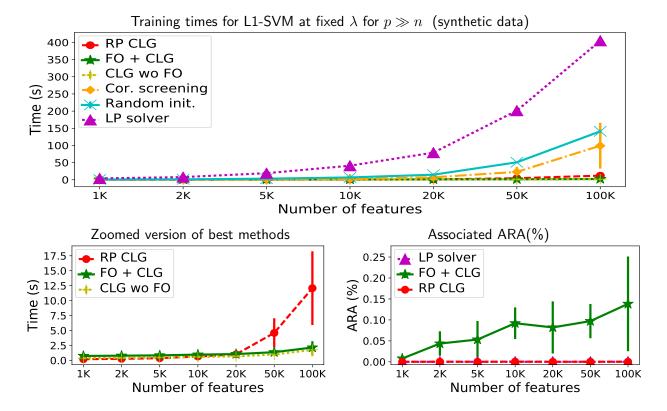


Figure 1: [Top panel] Run times (secs) for methods (a)-(e) for solving L1-SVM for a fixed λ with n=100 and varying p. [Bottom, left panel] Is a zoomed version of the top panel, where we show the best methods (a) and (b). [Bottom, right panel] presents the associated ARA values (in %), demonstrating that the solutions have similar accuracy levels. The basic version of the LP solver seems to be slow when compared to column generation, especially when p is large. The overall winner (both in terms of run time and solution quality) is F0+CLG (method (b)) — exhibiting up to a 100-times speedup over the LP-solver when $p=10^5$. The training time is equally split between first order and column generation algorithms. Method (b) solves a problem with $p=10^6$ in 10 seconds which more than 400 times faster than the LP solver (this is not shown in the figure so that it remains legible).

5.1.3 Computational results for n large and p small

Results for a fixed λ : We study the performance of the algorithms when $n \gg p$ — here, we anticipate constraint generation to be useful. We compare LP solver (Method (e)) described in Section 5.1.1 with our proposed method:

(f) "SFO+CNG": This is the constraint generation (CNG) method when initialized with a subsampling based first order (SFO) method heuristic (cf Section 4.4). For reference, we also report time taken by the column generation method alone (without SFO): "CNG wo SFO".

The subsampling based first order method is run sequentially across different subsamples: each subsample has a size $n_0 = 10p$, we use a tolerance $\mu_{\text{Tol}} = 10^{-1}$ and a maximum of $Q_{\text{max}} = n/n_0$ iterations. On each subsample, we consider a decreasing sequence of 5 values of τ (i.e., the smoothing parameter for the hinge-loss) with common ratio 0.7. The sample indices with nonzero

Training times for L1-SVM at fixed λ for $p \gg n$ (real datasets)

			FO+CI	LG (b)	LP solver (e)
Dataset	n	p	Time (s)	ARA (%)	Time (s)
Leukemia	72	7,129	0.59(0.02)	8.8×10^{-2}	20.94(1.08)
Lung cancer	181	12,533	1.45(0.02)	6.3×10^{-5}	87.92(2.75)
Ovarian	253	15,155	2.59(0.01)	1.2×10^{-2}	146.52(1.49)
Radsens	58	$12,\!625$	0.43(0.01)	1.6×10^{-1}	29.20(0.13)

Table 2: Training times in secs — we show the mean and the standard deviation of these ten values (within parenthesis) from: "F0+CLG" (method (b)) and "LP solver" (method (e)) for $\lambda = 0.01\lambda_{\rm max}$ on four real datasets. Column generation initialized with approximate solutions obtained from the first order method is found to be much faster than LP solver, in terms of reaching a near-optimal solution. The ARA for LP solver was approximately 0 for all instances.

hinge-loss are used to initialize constraint generation for a tolerance $\epsilon = 10^{-2}$.

To illustrate the behavior of our algorithms as a function of n, we take several values of n up to 50,000 with $p = 100, k_0 = 10$, $\rho = 0.1$ and set $\lambda = 0.01\lambda_{\text{max}}$. Figure 2 shows the results — we average the results over 10 simulations.

Training times for L1-SVM at fixed λ for $n \gg p$ on synthetic datasets

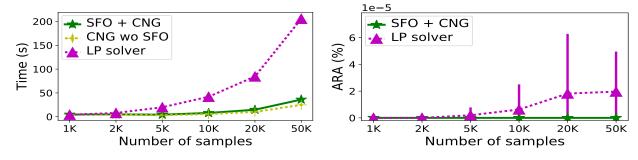


Figure 2: [Left] Training times for methods (e) and (f). [Right] Associated ARA(%). We set p = 100 and vary $n \in [1000, 50000]$. LP solver (e) is outperformed by our proposed Method (f) based on constraint generation. The quality of solutions in terms of objective values (see right panel plot for ARA) for both methods are similar. The overall training time of our method is dominated by the constraint generation algorithm (with little contribution from the initialization scheme).

5.1.4 Computational results for n large, p large

Results for a fixed λ on synthetic datasets: Sections 5.1.1–5.1.3 have shown the good performance of our column and constraint generation algorithms initialized with our first order heuristics. Figure 3 assesses the gain in performance when combining column and constraint generation methods to address instances with large values of n and p — this method is:

(g) "SF0+CL-CNG": We run the column-and-constraint generation algorithm (Algorithm 4 in Section 2.3.2) initialized by the subsampling heuristic presented in Section 4.4 (the method for large n and large p). For reference, we report CL-CNG wo SF0, which is the time taken to

perform the column-and-constraint generation algorithm excluding the initialization step.

Algorithm 4 is applied with $\epsilon = 0.01$. For the subsampling heuristic (cf Section 4.4, the method for large n and large p), we use the same parameter settings as in Section 5.1.3. Once the average estimate was obtained, we took the top 200 highest coefficients (in terms of absolute value) to initialize the set of columns for column generation.

We compare SF0+CL-CNG with Methods (a) and (b), which were the best-performing methods of Section 5.1.1. We set the regularization parameter for the L1-SVM problem to $\lambda = 0.001\lambda_{\text{max}}$. Figure 3 compares the averaged training times and ARA over 10 repetitions for n = 5000, $k_0 = 10$, $\rho = 0.1$ and a range of p values. For reference, we provide an account of the (average) number of columns and constraints that are active at an optimal solution: for p = 20K, the number of active features were ≈ 270 ; for p = 50K, the number of active features were ≈ 290 ; and for p = 100K, the number of active features were roughly similar.

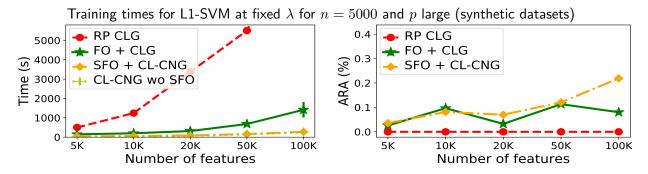


Figure 3: [Left] Training times for Methods (a), (b) and (g) for n = 5000. [Right] Associated ARA (%). The regularization path algorithm cannot handle large values of n and p, and basically explodes for $p = 10^5$. Our hybrid method (g) deals with a small subset of constraints as well as columns: hence, it significantly improves over column generation (alone), while returning an estimator with a similar objective value.

Results for a fixed λ on real datasets: Finally, we assess the quality of our hybrid column-and-constraint generation method (g) on two large real datasets¹⁴, when compared to the LP solver (e). All methods are run as explained above, except that we set $\mu_{\text{Tol}} = 0.5$. Because the real datasets are sparse, we use sparse matrices to deal with sparse matrix/vector multiplications—the sparsity in the LP model coefficient matrices are exploited by Gurobi's LP solvers. In this example, we do not standardize the covariates. We consider $\lambda = 0.05\lambda_{\text{max}}$ and initialize the set \mathcal{I} by retaining the 200 coefficients with largest magnitude, as available from the first order method. We average our results across 10 repetitions. The results are presented in Table 3.

¹⁴The datasets are from the UCI repository webpage https://archive.ics.uci.edu/ml/datasets.html?format=&task=cla&att=&area=&numAtt=&numIns=&type=&sort=nameUp&view=table.

Training times for L1-SVM on sparse real datasets at fixed λ for n large, p large

			SFO+CL-CNG	CL-CNG wo SFO	LP solver
Dataset	n	p	Time (s)	Time (s)	Time
rcv1 real-sim	,	,	188.67(2.88) 712.76(14.03)	75.41(2.91) 444.64(9.63)	>3 hrs >3 hrs

Table 3: Training times of our best Method (g) and the LP solver (e) on large sparse real datasets. The LP solver takes longer than 3 hours to converge, whereas our proposed methods converge substantially faster. Times taken by the first order methods is approximately half the total computation time of SFO+CL-CNG.

5.1.5 Comparison with a specialized state-of-the-art solver

The above experiments illustrate the gains in using the column/constraint generation framework over solving the full LP model. We compare our best methods — that is F0+CLG method (b) (for $p \gg n$ regime), SF0+CNG method (f) (for $n \gg p$ regime) versus a state-of-the-art algorithm PSM [26] which is a parametric simplex based solver (note that [26] demonstrate that PSM outperforms existing algorithms for solving the L1-SVM problem). In our experience, PSM is unable to handle instances where both n and p are large. As far as we can tell, PSM does not exploit column/constraint generation techniques. We use the solver made available by [26] with default parameter-settings. For this example (cf results in Table 4) computations for all methods were carried out on a MacBook with processor 1.6 GHz Intel Core i5, 8GB of RAM.

Training times for our L1-SVM methods versus state-of-the-art for $p\gg n$ and $n\gg p$ (synthetic datasets)

		Best Cutting Plane method			PSM		
n	p	Method	Time (s)	ARA (%)	Time (s)	ARA (%)	
100 100	10K 20K	FO+CLG FO+CLG	$1.53(0.22) \\ 1.89(0.28)$	0.07 (0.05) 0.08(0.04)	19.84(4.93) 50.78 (8.89)	$0.02(0.01) \\ 0.02(0.01)$	
1K 2K	100 100	SFO+CNG SFO+CNG	$3.27(0.43) \\ 3.32(1.03)$	0.00(0.00) 0.00(0.00)	49.85(7.91) 488.21(40.14)	$5.50(0.96) \\ 2.25(0.37)$	

Table 4: Training times of our best cutting plane methods compared to the state-of-the-art PSM solver on synthetic datasets with different values of n and p. PSM appears to be significantly slower than our methods (up to 150-times slower)—it also has a high variance in run time. When p is large, it leads to a small gain in accuracy, at the cost of a steep increase in training time. We note that PSM is unable to handle problems where both n, p are large, hence we do not report such examples.

5.2 Computational results for Group-SVM problem

We now study the performance of the column generation algorithm presented in Section 2.4 for the Group-SVM Problem (3).

Data Generation: Here, the covariates are drawn from a multivariate Gaussian with covariance Σ .

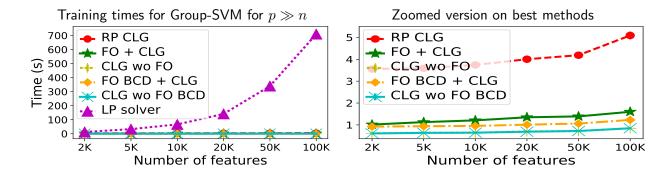


Figure 4: [Left panel] Comparison of training times for 4 methods (i)-(iv) for the Group-SVM Problem for $p \in [2 \times 10^3, 10^5]$ and n = 100. LP solver is much slower when compared to the other three methods. Our proposed column generation method initialized by a first order method, exhibits more than a 700-fold improvement in run time over LP solver. [Right panel] zooms in on the 3 best methods (i)-(iii). Accelerated gradient descent (ii) and block CD (iii) are both found to outperform the regularization path algorithm (i). Block CD (F0 BCD+CLG) appears as the overall winner: the run time improvement compared to F0+CLG is due to the CD method being faster than the gradient descent method.

The p covariates are divided into G groups each of the same size p_G . Within each group, covariates have pairwise correlation of ρ , and covariates are uncorrelated across groups (all variances are equal). Half of the samples are from the +1 class with (population) mean $\mu_+ = (\mathbf{1}_{p_G}, \dots, \mathbf{1}_{p_G}, \mathbf{0}_{p_G}, \dots, \mathbf{0}_{p_G})$; the remaining samples from class -1 have (population) mean $\mu_- = -\mu_+$.

Competing methods: We compare four algorithms for the Group-SVM Problem (3) with a fixed value of λ . We set $\lambda = 0.1\lambda_{\rm max}$ (with $\lambda_{\rm max}$ defined in (18)). The methods we compare are: (i): "RP CLG": This uses column generation with warm-start continuation across a grid of 6 equispaced regularization parameter values in $[\lambda_{\rm max}/2,\lambda]$ (cf Section 2.2.1). The column generation method uses a tolerance threshold of $\epsilon = 0.01$ (at every regularization parameter value). (ii): "F0+CLG": This applies column generation after initialization with a first order method (cf Section 4.3). We use a smoothing parameter $\tau = 0.2$ (for the hinge-loss) and use an accelerated gradient method restricted to the top n groups obtained via correlation screening (cf Section 4.4.1). We use a tolerance of $\epsilon = 0.01$ in the column generation method. We also report the run times of the column generation algorithms by excluding the times taken by the first order method by "CLG wo F0". (iii) "F0 BCD+CLG": This method is similar to (ii), except that we use a block CD method presented in Section 4.3 in place of the accelerated gradient method. We also report the times of running of the column generation algorithms without the block CD method: "CLG wo F0 BCD". (iv) "LP solver": This solves the full LP model with Gurobi's LP solver (no column generation is used).

Figure 4 presents the results — in this example, n = 100, $k_0 = 10$, $\rho = 0.1$ and each group has size $p_G = 10$. We average the results over 10 simulations. The ARA values for all the four methods are very similar, and are hence not reported.

5.3 Computational results for Slope-SVM

We present the computational performance of the column-and-constraint generation methods presented in Section 3.3 for the Slope-SVM Problem (4). The synthetic data is simulated as in Section 5.1.1 — we set n = 100, $k_0 = 10$, $\rho = 0.1$.

Comparison when λ_i s are not all distinct: We are not aware of any publicly available specialized implementations for the Slope SVM problem. We use the CVXPY modeling framework to model (See Section A.2) the Slope SVM problem and solve it using state-of-the art solvers like Ecos and Gurobi. We first consider a special instance of the Slope penalty (20) that corresponds to the coefficients $\lambda_i = 2\tilde{\lambda}$ for $i \leq k_0$ and $\lambda_i = \tilde{\lambda}$ for $i > k_0$; where $\tilde{\lambda} = 0.01\lambda_{\text{max}}$. We solve the resulting problem with both the Ecos and Gurobi solvers, denoted by "CVXPY Ecos" and "CVXPY Gurobi" (respectively). We compare them with our proposed column-and-constraint generation algorithm, referred to as "F0+CL-CNG". For our method, we first run the first order algorithm presented in Section 4.3 (for $\tau = 0.2$) restricted to the 10n columns with highest absolute correlations with the response (the remaining coefficients are all set to zero). The column-and-constraint generation algorithm (cf Section 3.3) uses a tolerance level of $\epsilon = 0.01$. We limit the number of columns to be added to 10 at each iteration. For reference we also display the run time of the algorithm by excluding the time taken by the initialization step and this is referred to as "CL-CNG wo F0". The results are presented in Table 5—we perform 10 replications for each of the methods.

Slope-SVM: our proposed methods versus CVXPY on synthetics datasets $(p \gg n = 100)$

	F0+0	CL-CNG	CL-CNG wo FO	CVXPY	Ecos	CVXPY	Gurobi
p	Time (s)	ARA (%)	Time (s)	Time (s)	ARA (%)	Time (s)	ARA (%)
10k	1.3(0.2)	4.8×10^{-2}	0.9(0.2)	38.6(3.4)	7.0×10^{-4}	58.4(0.6)	0.0
20k	1.6(0.3)	6.4×10^{-2}	2.3(0.6)	108.3(4.8)	0.0	130.3(1.4)	0.0
50k	2.7(0.6)	7.7×10^{-2}	2.3(0.6)	311.4(6.1)	6.0×10^{-3}	358.1(3.0)	0.0
100k	4.4(1.3)	9.6×10^{-2}	4.0(1.3)	_	_	757.2(5.6)	0.0

Table 5: Training times and ARA of our column-and-constraint generation method for Slope-SVM versus CVXPY—we took $\lambda_i/\lambda_j = 2$ for all $i \in [k_0]$ and $j > k_0$ (as described in the text). When the number of features are in the order of tens of thousands, our proposed method enjoys nearly a 200-fold speedup in run time. A '-' symbol denotes that the corresponding algorithm did not converge.

Comparison when λ_i s are distinct: We consider a general sequence of λ -values: Following [5], we set $\lambda_j = \sqrt{\log(2p/j)}\tilde{\lambda}$ with $\tilde{\lambda} = 0.01\lambda_{\text{max}}$. We observed that CVXPY could not handle even small instances of this problem (as the λ_j s are distinct) — in particular, the Ecos solver crashed for n = 100, p = 200. We compare our proposed method with the first order method (cf. Section 4.3) adapted to the full smooth Slope-SVM problem with $\tau = 0.2$. Due to the high per iteration cost of the first order method (FO), we terminate the method after a few iterations (with the associated

ARA within parenthesis). Table 6 compares our methods — we use the same synthetic dataset as in the previous case and average the results over 10 replications (the first order method was run for one replication due to its long run time).

Slope-SVM: our proposed methods versus first order methods on synthetics datasets $(p \gg n)$

	FO+CL-	-CNG	CL-CNG wo FO	First orde	er method (FO)
p	Time (s)	ARA (%)	Time (s)	Time (s)	ARA (%)
10k	10.84(5.63)	0.0(0.0)	9.57(5.63)	512.49	30.99
20k	22.82(10.12)	0.0(0.0)	21.24(10.12)	1071.24	32.42
50k	71.79(58.22)	0.0(0.0)	70.50(58.22)	> 3600	32.22

Table 6: Training times and ARA of our column-and-constraint generation for Slope-SVM with coefficients $\lambda_j = \sqrt{\log(2p/j)}\tilde{\lambda}$ (as in the text) on synthetic datasets with large number of features. Our proposed approach outperforms the first order methods (by at least 50-times) when they are asked to obtain a high accuracy solution. Due to the steep cost in training the first order methods, we ran them for only one replication (hence, we do not have any standard errors).

A Appendix

A.1 Methodology for column and constraint generation

We briefly review the methodology of column generation [13, 6]. The basic idea is to start with a candidate set of columns and incrementally add additional columns into the model until some optimality conditions are met. Consider the *primal* LP problem (**P**) where, \bar{n}, \bar{p} are integers and $\mathbf{A} \in \mathbb{R}^{\bar{n} \times \bar{p}}$, $\mathbf{b} \in \mathbb{R}^{\bar{n}}$, $\mathbf{c} \in \mathbb{R}^{\bar{p}}$ are problem data. We assume that the optimal objective value of (**P**) is finite. The strong duality theorem (cf [6], Theorem 4.4) states that this optimum is equal to that of the *dual* problem (**D**):

$$(\mathbf{P}): \min_{\boldsymbol{\theta} \in \mathbb{R}^{\bar{p}}} \mathbf{c}^{T}\boldsymbol{\theta}$$

$$\text{s.t.} \quad \mathbf{A}\boldsymbol{\theta} \ge \mathbf{b}, \quad \boldsymbol{\theta} \ge \mathbf{0}$$

$$(\mathbf{D}): \max_{\mathbf{q} \in \mathbb{R}^{\bar{n}}} \mathbf{q}^{T}\mathbf{b}$$

$$\text{s.t.} \quad \mathbf{q}^{T}\mathbf{A} \le \mathbf{c} \quad \mathbf{q} \ge \mathbf{0}.$$

$$(48)$$

We assume the rows of matrix \mathbf{A} to be linearly independent. We let \mathbf{A}_i denote the *i*th column of \mathbf{A} and \mathbf{a}_i denote its *i*th row. A subset $\mathcal{B} = \{B(1), \dots, B(\bar{n})\} \subset \{1, \dots, \bar{p}\}$ is said to define a basis if the columns of the matrix $\mathbf{B} = \{\mathbf{A}_{B(1)}, \dots, \mathbf{A}_{B(\bar{n})}\} \in \mathbb{R}^{\bar{n} \times \bar{n}}$ are linearly independent. The associated solution to the primal problem is $\boldsymbol{\theta} = (\boldsymbol{\theta}_{\mathcal{B}}, \mathbf{0}_{\bar{p} - \bar{n}})$ with $\boldsymbol{\theta}_{\mathcal{B}} = \mathbf{B}^{-1}\mathbf{b}$ and is feasible if $\mathbf{B}^{-1}\mathbf{b} \geq 0$. To lower the reduced cost $\mathbf{c}^T\boldsymbol{\theta}$, we select a non basic variable θ_j , $j \notin \mathcal{B}$ and increase it to a nonnegative value μ in a direction $\mathbf{d} \in \mathbb{R}^{\bar{p}}$ such that $\boldsymbol{\theta} + \mu \mathbf{d}$ is feasible for (\mathbf{P}) . The change in the objective value of Problem (\mathbf{P}) is equal to $\mu \bar{c}_j$ where \bar{c}_j is the reduced cost of the variable θ_j , defined as $\bar{c}_j = c_j - \mathbf{c}_{\mathcal{B}}^T \mathbf{B}^{-1} \mathbf{A}_j$. The simplex algorithm solves problem (\mathbf{P}) by looking for a non basic variable with negative reduced cost at every iteration. If such a variable cannot be found,

the current solution is optimal and the algorithm terminates. The basis \mathcal{B} also determines a basic solution to the dual Problem (**D**) given by $\mathbf{q}^T = \mathbf{c}_{\mathcal{B}}^T \mathbf{B}^{-1}$: it is feasible if $\mathbf{q} \geq 0$ and if all the reduced costs are nonnegative, since $\bar{c}_j = c_j - \mathbf{q}^T \mathbf{A}_j, \forall j$.

Column generation: We first consider the case $\bar{p} \gg \bar{n}$. A solution to Problem (**P**) has at most \bar{n} nonzeros coefficients. We anticipate that most of the columns will never enter the basis, hence we can reduce the number of "working" variables. Starting from a subset of variables $\mathcal{J} \subset \{1, \ldots, \bar{p}\}$, we define the *restricted columns* problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^{|\mathcal{J}|}} \qquad \sum_{j \in \mathcal{J}} c_j \theta_j
\text{s.t.} \qquad \sum_{j \in \mathcal{J}} \mathbf{A}_j \theta_j \ge \mathbf{b}, \ \boldsymbol{\theta} \ge 0.$$
(49)

The common version of the column generation algorithm adds to the set \mathcal{J} a non basic variable with lowest negative reduced cost and solves the updated restricted columns problem. The algorithm terminates after a finite number of iterations when no such variable can be found [cf 6, Section 6.2].

Constraint generation: We now consider the case when $\bar{n} \gg \bar{p}$. Suppose that, at an optimal solution to (**P**), only a small fraction of the \bar{n} constraints $\mathbf{a}_i^T \boldsymbol{\theta} \geq b_i$ for $i \in [\bar{n}]$ are active or binding—an optimal solution can be obtained by considering only a small subset of the \bar{n} constraints. This inspires the use of a constraint generation algorithm: we start from a subset of constraints $\mathcal{I} \subset \{1,\ldots,\bar{n}\}$, solve the restricted problem with these constraints and add the most violated constraint (if any). This method can also be interpreted as column generation [6] on the dual Problem (**D**).

A.2 Another formulation for the Slope norm

Without loss of generality, we consider a vector $\alpha \geq \mathbf{0}$. Now note that for every $m \in [p]$, we can represent $\alpha_{(1)} + \ldots + \alpha_{(m)} \leq s_m$ as

$$\alpha_{(1)} + \ldots + \alpha_{(m)} \le s_m, \quad \alpha \ge \mathbf{0} \iff \begin{cases} \mathbf{0} \le \alpha \le \theta_m \mathbf{1} + \mathbf{v}_m \\ m\theta_m + \mathbf{1}^T \mathbf{v}_m \le s_m \end{cases}$$
 (50)

with variables $\alpha, \mathbf{v}_m \in \mathbb{R}^p, \theta_m \in \mathbb{R}$ and $\mathbf{1} \in \mathbb{R}^p$ being a vector of all ones. Note that the rhs formulation in (50) has O(p) variables and O(p)-constraints. Note that we can write:

$$\sum_{j=1}^{p} \lambda_j \alpha_{(j)} = \sum_{m=1}^{p} \tilde{\lambda}_m (\alpha_{(1)} + \ldots + \alpha_{(m)})$$

where, $\tilde{\lambda}_m = \lambda_m - \lambda_{m-1}$ for all $m \in \{1, ..., p\}$. Therefore, representing $\sum_{j=1}^p \lambda_j \alpha_{(j)} \leq \eta$ will require a representation (50) for m = 1, ..., p — this will lead to a formulation with $O(p^2)$ variables and $O(p^2)$ constraints, which can be quite large as soon as p becomes a few hundred.

References

- [1] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Convex optimization with sparsity-inducing norms. Optimization for Machine Learning, 5:19–53, 2011.
- [2] P. Balamurugan, A. Posinasetty, and S. Shevade. Admm for training sparse structural syms with augmented l1 regularizers. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, pages 684–692. SIAM, 2016.
- [3] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM journal on imaging sciences, 2(1):183–202, 2009.
- [4] S. R. Becker, E. J. Candès, and M. C. Grant. Templates for convex cone problems with applications to sparse signal recovery. *Mathematical programming computation*, 3(3):165, 2011.
- [5] P. C. Bellec, G. Lecué, A. B. Tsybakov, et al. Slope meets lasso: improved oracle bounds and optimality. The Annals of Statistics, 46(6B):3603–3642, 2018.
- [6] D. Bertsimas and J. N. Tsitsiklis. Introduction to linear optimization, volume 6. Athena Scientific Belmont, MA, 1997.
- [7] M. Bogdan, E. van den Berg, C. Sabatti, W. Su, and E. J. Candès. Slope adaptive variable selection via convex optimization. *The annals of applied statistics*, 9(3):1103, 2015.
- [8] L. Bottou. Large-scale machine learning with stochastic gradient descent. In Proceedings of COMP-STAT'2010, pages 177–186. Springer, 2010.
- [9] S. Boyd and L. Vandenberghe. Convex optimization. Cambridge university press, 2004.
- [10] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- [11] P. S. Bradley and O. L. Mangasarian. Feature selection via concave minimization and support vector machines. In *ICML*, volume 98, pages 82–90, 1998.
- [12] J. Desrosiers and M. E. Lübbecke. A primer in column generation. In Column generation, pages 1–32. Springer, 2005.
- [13] L. R. Ford Jr and D. R. Fulkerson. A suggested computation for maximal multi-commodity network flows. *Management Science*, 5(1):97–101, 1958.

- [14] V. Franc and S. Sonnenburg. Optimized cutting plane algorithm for support vector machines. In Proceedings of the 25th international conference on Machine learning, pages 320–327. ACM, 2008.
- [15] J. Friedman, T. Hastie, and R. Tibshirani. Regularized paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1), 2010.
- [16] T. Hastie, S. Rosset, R. Tibshirani, and J. Zhu. The entire regularization path for the support vector machine. *Journal of Machine Learning Research*, 5(Oct):1391–1415, 2004.
- [17] T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer New York, 2 edition, 2009.
- [18] C.-J. Hsieh, K.-W. Chang, C.-J. Lin, S. S. Keerthi, and S. Sundararajan. A dual coordinate descent method for large-scale linear sym. In *Proceedings of the 25th international conference on Machine* learning, pages 408–415. ACM, 2008.
- [19] J. Huang and T. Zhang. The benefit of group sparsity. The Annals of Statistics, 38(4):1978–2004, 2010.
- [20] T. Joachims. Training linear syms in linear time. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 217–226. ACM, 2006.
- [21] J. D. Lee, Q. Liu, Y. Sun, and J. E. Taylor. Communication-efficient sparse regression. Journal of Machine Learning Research, 18(5):1–30, 2017.
- [22] J.-J. Moreau. Dual convex functions and proximal points in a hilbert space. CR Acad. Sci. Paris Ser. At Math. volume = 255, pages = 2897–2899 year = 1962.
- [23] Y. Nesterov. Introductory Lectures on Convex Optimization: A Basic Course. Kluwer, Norwell, 2004.
- [24] Y. Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.
- [25] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140 (1):125–161, 2013.
- [26] H. Pang, H. Liu, R. J. Vanderbei, and T. Zhao. Parametric simplex method for sparse learning. In Advances in Neural Information Processing Systems, pages 188–197, 2017.
- [27] Z. Qin, K. Scheinberg, and D. Goldfarb. Efficient block-coordinate descent algorithms for the group lasso. *Mathematical Programming Computation*, 5(2):143–169, 2013.
- [28] T. Robertson and T. Robertson. Order restricted statistical inference. Wiley, New York., 1988.
- [29] S. Shalev-Shwartz, Y. Singer, and N. Srebro. Pegasos: Primal estimated sub-gradient solver for svm. In Proceedings of the 24th international conference on Machine learning, pages 807–814. ACM, 2007.

- [30] R. Tibshirani, J. Bien, J. Friedman, T. Hastie, N. Simon, J. Taylor, and R. J. Tibshirani. Strong rules for discarding predictors in lasso-type problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(2):245–266, 2012.
- [31] S. A. van de Geer. Empirical Processes in M-estimation, volume 6. Cambridge University Press, 2000.
- [32] E. van den Berg and M. P. Friedlander. SPGL1: A solver for large-scale sparse reconstruction, June 2007. http://www.cs.ubc.ca/labs/scl/spgl1.
- [33] E. van den Berg and M. P. Friedlander. Probing the pareto frontier for basis pursuit solutions. SIAM Journal on Scientific Computing, 31(2):890-912, 2008. doi: 10.1137/080714488. URL http://link.aip.org/link/?SCE/31/890.
- [34] V. Vapnik. The nature of statistical learning theory. Springer science & business media, 2013.
- [35] S. J. Wright. Coordinate descent algorithms. Mathematical Programming, 151(1):3-34, 2015.
- [36] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- [37] X. Zeng and M. Figueiredo. The ordered weighted 11 norm: Atomic formulation. *Dual Norm, and Projections. arXiv preprint.*