

The 2014–2015 Philip McCord Morse Lecture

**Statistics and Machine Learning
via a Modern Optimization Lens**

Dimitris Bertsimas

Boeing Professor of Operations Research and Statistics
Sloan School of Management and Operations Research Center,
Massachusetts Institute of Technology,
Cambridge, Massachusetts

November 11, 2014
INFORMS Annual Meeting
San Francisco, CA

Best Subset Selection via a Modern Optimization Lens

Dimitris Bertsimas

Sloan School of Management and Operations Research Center,
Massachusetts Institute of Technology,
Cambridge, Massachusetts

Angela King

Operations Research Center, Massachusetts Institute of Technology,
Cambridge, Massachusetts

Rahul Mazumder

Department of Statistics, Columbia University,
New York, New York

INFORMS Philip Morse Lecture for 2014–2015

Abstract

In the last 25 years (1990–2014), algorithmic advances in integer optimization combined with hardware improvements have resulted in an astonishing 200 billion factor speedup in solving mixed integer optimization (MIO) problems. We present a MIO approach for solving the classical best subset selection problem of choosing k out of p features in linear regression given n observations. We develop a discrete extension of modern first order continuous optimization methods to find high quality feasible solutions that we use as warm starts to a MIO solver that finds provably optimal solutions. The resulting algorithm (a) provides a solution with a guarantee on its suboptimality even if we terminate the algorithm early, (b) can accommodate side constraints on the coefficients of the linear regression, and (c) extends to finding best subset solutions for the least absolute deviation loss function. Using a wide variety of synthetic and real datasets, we demonstrate that our approach solves problems with n in the 1000s and p in the 100s in minutes to provable optimality, and finds near optimal solutions for n in the 100s and p in the 1000s in minutes. We also establish via numerical experiments that the MIO approach performs better than Lasso in terms of achieving sparse solutions with good predictive power.

1 Introduction

We consider the linear regression model with response vector $\mathbf{y}_{n \times 1}$, model matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p] \in \mathbb{R}^{n \times p}$, regression coefficients $\boldsymbol{\beta} \in \mathbb{R}^{p \times 1}$ and errors $\boldsymbol{\epsilon} \in \mathbb{R}^{n \times 1}$:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

In many important classical and modern statistical applications, it is desirable to obtain a parsimonious fit to the data by finding the best k -feature fit to the response \mathbf{y} . We will also assume that the columns of \mathbf{X} have been standardized to have zero means and unit ℓ_2 -norm. Especially in the high-dimensional regime with $p \gg n$, in order to conduct statistically meaningful inference, it is desirable to assume that the true regression coefficient $\boldsymbol{\beta}$ is sparse or may be well approximated by a sparse vector. Quite naturally, the last few decades have seen a flurry of activity in estimating sparse linear models with good explanatory power. Central to this statistical task lies the best subset Problem [35] with subset size k , which is given by the following optimization problem:

$$\min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_0 \leq k, \quad (1)$$

where the ℓ_0 (pseudo)norm of a vector $\boldsymbol{\beta}$ counts the number of nonzeros in $\boldsymbol{\beta}$ and is given by $\|\boldsymbol{\beta}\|_0 = \sum_{i=1}^p 1(\beta_i \neq 0)$, where $1(\cdot)$ denotes the indicator function. The cardinality constraint makes Problem (1) NP-hard [36]. Indeed, state-of-the-art algorithms to solve Problem (1), as implemented in popular statistical packages, like `leaps` in `R`, do not scale to problem sizes larger than $p = 30$. Due to this reason, it is not surprising that the best subset problem has been widely dismissed as being *intractable* by the greater statistical community.

In this paper we address Problem (1) using modern optimization methods, specifically mixed integer optimization (MIO) and a discrete extension of first order continuous optimization methods. Using a wide variety of synthetic and real datasets, we demonstrate that our approach solves problems with n in the 1000s and p in the 100s in minutes to provable optimality, and finds near optimal solutions for n in the 100s and p in the 1000s in minutes. To the best of our knowledge, this is the first time that MIO has been demonstrated to be a tractable solution method for Problem (1). We note that we use the term tractability not to mean the usual polynomial solvability for problems, but rather the ability to solve problems of realistic size in times that are appropriate for the applications we consider.

As there is a vast literature on the best subset problem, we next give a brief and selective overview of related approaches for the problem.

Brief Context and Background

To overcome the computational difficulties of the best subset problem, computationally tractable convex optimization based methods like Lasso [44, 15] have been proposed as a convex surrogate for Problem (1). For the linear regression problem, the Lagrangian form of Lasso solves

$$\min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1, \quad (2)$$

where the ℓ_1 penalty on $\boldsymbol{\beta}$, i.e., $\|\boldsymbol{\beta}\|_1 = \sum_i |\beta_i|$ shrinks the coefficients towards zero and naturally produces a sparse solution by setting many coefficients to be exactly zero. There has been a substantial amount of impressive work on Lasso [21, 13, 5, 49, 28, 53, 17, 31, 34, 47, 45] in terms of algorithms and understanding of its theoretical properties—see for example the excellent books or surveys [10, 30, 45] and the references therein.

Indeed, Lasso enjoys several attractive statistical properties and has drawn a significant amount of attention from the statistics community as well as other closely related fields. Under various conditions on the model matrix \mathbf{X} and $n, p, \boldsymbol{\beta}$ it can be shown that Lasso delivers a sparse model with good predictive performance [10, 30]. In order to perform exact variable selection, much stronger assumptions are required [10]. Sufficient conditions under which Lasso gives a sparse model with good predictive performance are the restricted eigenvalue conditions and compatibility conditions [10]. These involve statements about the range of the spectrum of sub-matrices of \mathbf{X} and are difficult to verify, for a given data-matrix \mathbf{X} .

An important reason behind the popularity of Lasso is its computational feasibility and scalability to practical sized problems. Problem (2) is a convex quadratic optimization problem and there are several efficient solvers for it, see for example [39, 21, 25].

In spite of its favorable statistical properties, Lasso has several shortcomings. In the presence of noise and correlated variables, in order to deliver a model with good predictive accuracy, Lasso brings in a large number of nonzero coefficients (all of which are shrunk towards zero) including noise variables. Lasso leads to biased regression coefficient estimates, since the ℓ_1 -norm penalizes both large and small coefficients uniformly. In contrast, if the best subset selection procedure decides to include a variable in the model, it brings it in without any shrinkage thereby draining the effect of its correlated surrogates. Upon increasing the degree of regularization, Lasso sets more coefficients to zero, but in the process ends up leaving out true predictors from the active set. Thus, as soon as certain sufficient regularity conditions on the data are violated, Lasso becomes suboptimal as (a) a variable selector and (b) in terms of delivering a model with good predictive performance.

The shortcomings of Lasso are also known in the statistical literature. In fact, there is a significant gap between what can be achieved via best subset selection and Lasso: this is supported by empirical (for small problem sizes, i.e., $p \leq 30$) and theoretical evidence, see for example, [41, 52, 33, 27, 50, 43] and the references therein.

To address the shortcomings, non-convex penalized regression is often used to “bridge” the gap between the convex ℓ_1 penalty and the combinatorial ℓ_0 penalty [33, 23, 22, 48, 49, 24, 54, 55, 51, 11]. Written in Lagrangian form, this gives rise to continuous non-convex optimization problems of the form:

$$\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \sum_i p(|\beta_i|; \gamma; \lambda), \quad (3)$$

where $p(|\beta|; \gamma; \lambda)$ is a non-convex function in β with λ and γ denoting the degree of regularization and non-convexity, respectively. Typical examples

of non-convex penalties include the minimax concave penalty (MCP), the smoothly clipped absolute deviation (SCAD), and ℓ_γ penalties (see for example, [23, 33, 55, 22]). There is strong statistical evidence indicating the usefulness of estimators obtained as minimizers of non-convex penalized problems (3) over Lasso see for example [50, 32, 48].

Problem (3) mainly leads to a family of continuous and non-convex optimization problems. Various effective nonlinear optimization based methods (see for example [55, 22, 11, 32, 48, 33] and the references therein) have been proposed in the literature to obtain good local minimizers to Problem (3). In particular [33] proposes Sparsenet, a coordinate-descent procedure to trace out a surface of local minimizers for Problem (3) for the MCP penalty using effective warm start procedures. None of the existing approaches for solving Problem (3), however, come with guarantees of how close the solutions are to the global minimum of Problem (3).

The Lagrangian version of (1) given by

$$\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{i=1}^p 1(\beta_i \neq 0), \quad (4)$$

may be seen as a special case of (3). Note that, due to non-convexity, problems (4) and (1) are *not* equivalent. Problem (1) allows one to control the exact level of sparsity via the choice of k , unlike (4) where there is no clear correspondence between λ and k . Problem (4) is a discrete optimization problem unlike continuous optimization problems (3) arising from continuous non-convex penalties.

Insightful statistical properties of Problem (4) have been explored from a theoretical viewpoint in [50, 27, 28, 43]. [43] points out that (1) is preferable over (4) in terms of superior statistical properties of the resulting estimator. None of the aforementioned papers, however, discuss methods to obtain provably optimal solutions to problems (4) or (1), and to the best of our knowledge, computing optimal solutions to problems (4) and (1) is deemed as intractable.

Our Approach In this paper, we propose a novel framework via which the best subset selection problem can be solved to optimality or near optimality in problems of practical interest within a reasonable time frame. At the core of our proposal is a computationally tractable framework that brings to bear the power of modern discrete optimization methods: discrete first order methods motivated by first order methods in convex optimization [40] and mixed integer optimization (MIO), see [4]. We do not guarantee polynomial time solution times as these do not exist for the best subset problem unless P=NP. Rather, our view of computational tractability is the ability of a method to solve problems of practical interest in times that are appropriate for the application addressed. An advantage of our approach is that it adapts to variants of the best subset regression problem of the form:

$$\begin{aligned} \min_{\boldsymbol{\beta}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_q^q \\ \text{s.t.} \quad & \|\boldsymbol{\beta}\|_0 \leq k \\ & \mathbf{A}\boldsymbol{\beta} \leq \mathbf{b}, \end{aligned}$$

where $\mathbf{A}\boldsymbol{\beta} \leq \mathbf{b}$ represents polyhedral constraints and $q \in \{1, 2\}$ refers to a least absolute deviation or the least squares loss function on the residuals $\mathbf{r} := \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$.

Existing approaches in the Mathematical Optimization Literature

In a seminal paper [26], the authors describe a leaps and bounds procedure for computing global solutions to Problem (1) (for the classical $n > p$ case) which can be achieved with computational effort significantly less than complete enumeration. leaps, a state-of-the-art R package uses this principle to perform best subset selection for problems with $n > p$ and $p \leq 30$. [3] proposed a tailored branch-and-bound scheme that can be applied to Problem (1) using ideas from [26] and techniques in quadratic optimization, extending and enhancing the proposal of [6]. The proposal of [3] concentrates on obtaining high quality upper bounds for Problem (1) and is less scalable than the methods presented in this paper.

Contributions We summarize our contributions in this paper below:

1. We use MIO to find a provably optimal solution for the best subset problem. Our approach has the appealing characteristic that if we terminate the algorithm early, we obtain a solution with a guarantee on its suboptimality. Furthermore, our framework can accommodate side constraints on $\boldsymbol{\beta}$ and also extends to finding best subset solutions for the least absolute deviation loss function.
2. We introduce a general framework of solution methods based on a discrete extension of modern first order continuous optimization methods that provide near-optimal solutions for the best subset problem. The MIO algorithm significantly benefits from solutions obtained by the first order methods and problem specific information that can be computed in a data-driven fashion.
3. We report computational results with both synthetic and real-world datasets that show that our proposed framework can deliver provably optimal solutions for problems of size n in the 1000s and p in the 100s in minutes. For high-dimensional problems with $n \in \{50, 100\}$ and $p \in \{1000, 2000\}$, with the aid of warm starts and further problem-specific information, our approach finds near optimal solutions in minutes but takes hours to prove optimality.
4. We investigate the statistical properties of best subset selection procedures for practical problem sizes, which to the best of our knowledge, have remained largely unexplored to date. We demonstrate the favorable predictive performance and sparsity-inducing properties of the best subset selection procedures over Lasso, Sparsenet and stepwise regression in a wide variety of real and synthetic examples for both the least squares and absolute deviation loss functions.

The structure of the paper is as follows. In Section 2, we present a brief overview of MIO, including a summary of the computational advances it has enjoyed in the last twenty-five years. We present the proposed MIO formulations for the best subset problem as well as some connections with the compressed sensing literature for estimating parameters and providing lower bounds for the MIO formulations that improve their computational performance. In Section 3, we develop a discrete extension of first order methods in convex optimization to obtain near optimal solutions for the best subset problem and establish its convergence properties, a method that may be of independent interest. In Section 4, we perform a variety of computational tests on synthetic and real datasets to assess the algorithmic and statistical performances of our approach for the least squares loss function for both the classical overdetermined case $n > p$, and the high-dimensional case $p \gg n$. In Section 5, we report computational results for the least absolute deviation loss function. In Section 6, we include our concluding remarks.

2 Mixed Integer Optimization Formulations

In this section, we present a brief overview of MIO, including the simply astonishing advances it has enjoyed in the last twenty-five years. We then present the proposed MIO formulations for the best subset problem as well as some connections with the compressed sensing literature for estimating parameters and providing lower bounds for the MIO formulations that improve their computational performance.

2.1 Brief Background on MIO

The general form of a Mixed Integer Quadratic Optimization (MIQO) problem is as follows:

$$\begin{aligned} \min \quad & \boldsymbol{\alpha}^T \mathbf{Q} \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{a} \\ \text{s.t.} \quad & \mathbf{A} \boldsymbol{\alpha} \leq \mathbf{b} \\ & \alpha_i \in \{0, 1\}, \quad \forall i \in \mathcal{I} \\ & \alpha_j \in \mathbb{R}_+, \quad \forall j \notin \mathcal{I}, \end{aligned}$$

where $\mathbf{a} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{k \times m}$, $\mathbf{b} \in \mathbb{R}^k$ and $\mathbf{Q} \in \mathbb{R}^{m \times m}$ (positive semidefinite) are the given parameters of the problem; \mathbb{R}_+ denotes the non-negative reals, the symbol \leq denotes element-wise inequalities and we optimize over $\boldsymbol{\alpha} \in \mathbb{R}^m$ containing both discrete ($\alpha_i, i \in \mathcal{I}$) and continuous ($\alpha_i, i \notin \mathcal{I}$) variables, with $\mathcal{I} \subset \{1, \dots, m\}$. For background on MIO see [4]. Subclasses of MIQO problems include convex quadratic optimization problems ($\mathcal{I} = \emptyset$), mixed integer ($\mathbf{Q} = \mathbf{0}_{m \times m}$) and linear optimization problems ($\mathcal{I} = \emptyset, \mathbf{Q} = \mathbf{0}_{m \times m}$). Modern integer optimization solvers such as Gurobi and CPLEX are able to tackle MIQO problems.

In the last twenty-five years (1991-2014) the computational power of MIO solvers has increased at an astonishing rate. In [7], to measure the speedup of MIO solvers, the same set of MIO problems were tested on the same computers

using twelve consecutive versions of CPLEX and version-on-version speedups were reported. The versions tested ranged from CPLEX 1.2, released in 1991 to CPLEX 11, released in 2007. Each version released in these years produced a speed improvement on the previous version, leading to a total speedup factor of more than 29,000 between the first and last version tested (see [7], [37] for details). Gurobi 1.0, a MIO solver which was first released in 2009, was measured to have similar performance to CPLEX 11. Version-on-version speed comparisons of successive Gurobi releases have shown a speedup factor of more than 20 between Gurobi 5.5, released in 2013, and Gurobi 1.0 ([7], [37]). The combined machine-independent speedup factor in MIO solvers between 1991 and 2013 is 580,000. This impressive speedup factor is due to incorporating both theoretical and practical advances into MIO solvers. Cutting plane theory, disjunctive programming for branching rules, improved heuristic methods, techniques for preprocessing MIOs, using linear optimization as a black box to be called by MIO solvers, and improved linear optimization methods have all contributed greatly to the speed improvements in MIO solvers [7].

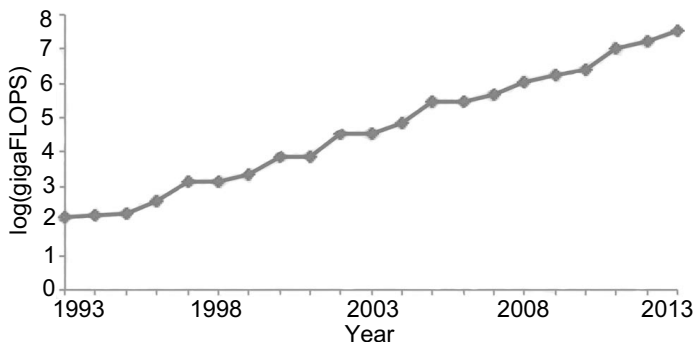


Figure 1: Log of Peak Supercomputer Speed from 1993–2013.

In addition, the past twenty years have also brought dramatic improvements in hardware. Figure 1 shows the exponentially increasing speed of supercomputers over the past twenty years, measured in billion floating point operations per second [1]. The hardware speedup from 1993 to 2013 is approximately $10^{5.5} \sim 320,000$. When both hardware and software improvements are considered, the overall speedup¹ is approximately 200 billion! MIO solvers provide both feasible solutions as well as lower bounds to the optimal value. As the MIO solver progresses towards the optimal solution, the lower bounds improve and provide an increasingly better guarantee of suboptimality, which is especially useful if the MIO solver is stopped before reaching the global optimum. In contrast, heuristic methods do not provide such a certificate of suboptimality.

¹Note that the speedup factors cited here refer to mixed integer linear optimization problems, not MIQO problems. The speedup factors for MIQO problems are similar.

The belief that MIO approaches to problems in statistics are not practically relevant was formed in the 1970s and 1980s and it was at the time justified. Given the astonishing speedup of MIO solvers and computer hardware in the last twenty-five years, the mindset of MIO as theoretically elegant but practically irrelevant is no longer justified. In this paper, we provide empirical evidence of this fact in the context of the best subset selection problem.

2.2 MIO Formulations for the Best Subset Selection Problem

We first present a simple reformulation to Problem (1) as a MIO (in fact a MIQO) problem:

$$\begin{aligned}
 Z_1 = \min_{\boldsymbol{\beta}, \mathbf{z}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \\
 \text{s.t.} \quad & -\mathcal{M}_U z_i \leq \beta_i \leq \mathcal{M}_U z_i, i = 1, \dots, p \\
 & z_i \in \{0, 1\}, i = 1, \dots, p \\
 & \sum_{i=1}^p z_i \leq k,
 \end{aligned} \tag{5}$$

where $\mathbf{z} \in \{0, 1\}^p$ is a binary variable and \mathcal{M}_U is a constant such that if $\widehat{\boldsymbol{\beta}}$ is a minimizer of Problem (5), then $\mathcal{M}_U \geq \|\widehat{\boldsymbol{\beta}}\|_\infty$. If $z_i = 1$, then $|\beta_i| \leq \mathcal{M}_U$ and if $z_i = 0$, then $\beta_i = 0$. Thus, $\sum_{i=1}^p z_i$ is an indicator of the number of zeros in $\boldsymbol{\beta}$.

Provided that \mathcal{M}_U is chosen to be sufficiently large with $\mathcal{M}_U \geq \|\widehat{\boldsymbol{\beta}}\|_\infty$, a solution to Problem (5) will be a solution to Problem (1). Of course, \mathcal{M}_U is not known a priori, and a small value of \mathcal{M}_U may lead to a solution different from (1). The choice of \mathcal{M}_U affects the strength of the formulation and is critical for obtaining solutions quickly in practice. In Section 2.3 we describe how to find appropriate values for \mathcal{M}_U .

Formulation (5) leads to interesting insights, especially via the structure of the convex hull of its constraints, as illustrated next :

$$\begin{aligned}
 & \text{Conv} \left(\left\{ \boldsymbol{\beta} : |\beta_i| \leq \mathcal{M}_U z_i, z_i \in \{0, 1\}, i = 1, \dots, p, \sum_{i=1}^p z_i \leq k \right\} \right) \\
 &= \{ \boldsymbol{\beta} : \|\boldsymbol{\beta}\|_\infty \leq \mathcal{M}_U, \|\boldsymbol{\beta}\|_1 \leq \mathcal{M}_U k \} \\
 &\subseteq \{ \boldsymbol{\beta} : \|\boldsymbol{\beta}\|_1 \leq \mathcal{M}_U k \}.
 \end{aligned}$$

Thus, the minimum of Problem (5) is lower-bounded by the optimum objective value of both the following convex optimization problems:

$$Z_2 := \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_\infty \leq \mathcal{M}_U, \|\boldsymbol{\beta}\|_1 \leq \mathcal{M}_U k \tag{6}$$

$$Z_3 := \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_1 \leq \mathcal{M}_U k, \tag{7}$$

where (7) is the familiar Lasso in constrained form. This is a weaker relaxation than formulation (6), which in addition to the ℓ_1 constraint on β , has box-constraints controlling the values of the β_i 's. It is easy to see that the following ordering exists:

$$Z_3 \leq Z_2 \leq Z_1,$$

with the inequalities being strict in most instances.

In terms of approximating the optimal solution to Problem (5), the MIO solver begins by first solving a continuous relaxation of Problem (5). The Lasso formulation (7) is weaker than this root node relaxation. Additionally, MIO is typically able to significantly improve the quality of the root node solution as the MIO solver progresses toward the optimal solution.

To motivate the reader we provide an example of the evolution (see Figure 2) of the MIO formulation (8) for the Diabetes dataset [21], with $n = 350, p = 64$ (for further details on the dataset see Section 4).

Since formulation (5) is sensitive to the choice of \mathcal{M}_U , we consider an alternative MIO formulation based on Specially Ordered Sets [4] as described next.

Formulations via Specially Ordered Sets Any feasible solution to formulation (5) will have $(1 - z_i)\beta_i = 0$ for every $i \in \{1, \dots, p\}$. This constraint can be modeled via integer optimization using Specially Ordered Sets of Type 1 [4] (SOS-1). In an SOS-1 constraint, at most one variable in the set can take a nonzero value, that is

$$(1 - z_i)\beta_i = 0 \iff (\beta_i, 1 - z_i) : \text{SOS-1},$$

for every $i = 1, \dots, p$. This leads to the following formulation of (1):

$$\begin{aligned} \min_{\beta, \mathbf{z}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \\ \text{s.t.} \quad & (\beta_i, 1 - z_i) : \text{SOS-1}, \quad i = 1, \dots, p \\ & z_i \in \{0, 1\}, \quad i = 1, \dots, p \\ & \sum_{i=1}^p z_i \leq k, \end{aligned} \tag{8}$$

The objective function in Problem (8) is a convex quadratic function in the continuous variable β , which can be formulated explicitly as:

$$\begin{aligned} \min_{\beta, \mathbf{z}} \quad & \frac{1}{2} \beta^T \mathbf{X}^T \mathbf{X} \beta - \langle \mathbf{X}'\mathbf{y}, \beta \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 \\ \text{s.t.} \quad & (\beta_i, 1 - z_i) : \text{SOS-1}, \quad i = 1, \dots, p \\ & z_i \in \{0, 1\}, \quad i = 1, \dots, p \\ & \sum_{i=1}^p z_i \leq k \\ & -\mathcal{M}_U \leq \beta_i \leq \mathcal{M}_U, \quad i = 1, \dots, p \\ & \|\beta\|_1 \leq \mathcal{M}_\ell. \end{aligned} \tag{9}$$

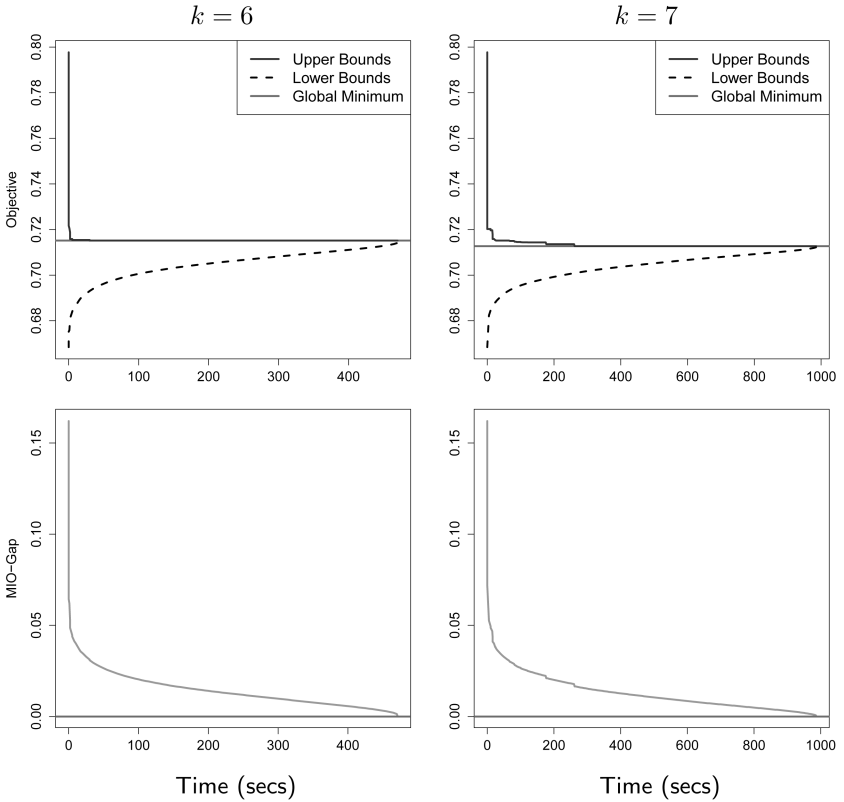


Figure 2: The typical evolution of the MIO formulation (8) for the diabetes dataset with $n = 350, p = 64$ with $k = 6$ (left panel) and $k = 7$ (right panel). The top panel shows the evolution of upper and lower bounds with time. The lower panel shows the evolution of the corresponding MIO gap with time. Optimal solutions for both the problems are found in a few seconds in both examples, but it takes 10-20 minutes to certify optimality via the lower bounds. Note that the time taken for the MIO to certify convergence to the global optimum increases with increasing k .

We also provide problem-dependent constants \mathcal{M}_U and $\mathcal{M}_\ell \in [0, \infty]$. \mathcal{M}_U provides an upper bound on the absolute value of the regression coefficients and \mathcal{M}_ℓ provides an upper bound on the ℓ_1 -norm of β . Adding these bounds typically leads to improved performance of the MIO. In Section 2.3, we describe an approach to compute these parameters from the data.

We also consider another formulation for (9):

$$\begin{aligned}
\min_{\boldsymbol{\beta}, \mathbf{z}, \boldsymbol{\zeta}} \quad & \frac{1}{2} \boldsymbol{\zeta}^T \boldsymbol{\zeta} - \langle \mathbf{X}' \mathbf{y}, \boldsymbol{\beta} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 \\
\text{s.t.} \quad & \boldsymbol{\zeta} = \mathbf{X} \boldsymbol{\beta} \\
& (\beta_i, 1 - z_i) : \text{SOS-1}, \quad i = 1, \dots, p \\
& z_i \in \{0, 1\}, \quad i = 1, \dots, p \\
& \sum_{i=1}^p z_i \leq k \\
& -\mathcal{M}_U \leq \beta_i \leq \mathcal{M}_U, \quad i = 1, \dots, p \\
& \|\boldsymbol{\beta}\|_1 \leq \mathcal{M}_\ell \\
& -\mathcal{M}_U^\zeta \leq \zeta_i \leq \mathcal{M}_U^\zeta, \quad i = 1, \dots, n \\
& \|\boldsymbol{\zeta}\|_1 \leq \mathcal{M}_\ell^\zeta,
\end{aligned} \tag{10}$$

where the optimization variables are $\boldsymbol{\beta} \in \mathbb{R}^p, \boldsymbol{\zeta} \in \mathbb{R}^n, \mathbf{z} \in \{0, 1\}^p$ and $\mathcal{M}_U, \mathcal{M}_\ell, \mathcal{M}_U^\zeta, \mathcal{M}_\ell^\zeta \in [0, \infty]$ are problem specific parameters. Note that the objective function in formulation (10) involves a quadratic form in n variables and a linear function in p variables. Problem (10) is equivalent to the following variant of the best subset problem:

$$\begin{aligned}
\min_{\boldsymbol{\beta}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|_2^2 \\
\text{s.t.} \quad & \|\boldsymbol{\beta}\|_\infty \leq \mathcal{M}_U, \|\boldsymbol{\beta}\|_1 \leq \mathcal{M}_\ell \\
& \|\mathbf{X} \boldsymbol{\beta}\|_\infty \leq \mathcal{M}_U^\zeta, \|\mathbf{X} \boldsymbol{\beta}\|_1 \leq \mathcal{M}_\ell^\zeta.
\end{aligned} \tag{11}$$

Formulations (9) and (10) differ in the size of the quadratic forms that are involved. The current state-of-the-art MIO solvers are better-equipped to handle mixed integer linear optimization problems than MIQO problems. Formulation (9) has fewer variables but a quadratic form in p variables—we find this formulation more useful in the $n > p$ regime, with p in the 100s. Formulation (10) on the other hand has more variables, but involves a quadratic form in n variables—this formulation is more useful for high-dimensional problems $p \gg n$, with n in the 100s and p in the 1000s.

The bounds on $\boldsymbol{\beta}$ and $\boldsymbol{\zeta}$ are not required, but if these constraints are provided, they improve the strength of the MIO formulation. We next show how these bounds can be computed from given data.

2.3 Specification of Parameters

In this section, we obtain estimates for the quantities $\mathcal{M}_U, \mathcal{M}_\ell, \mathcal{M}_U^\zeta, \mathcal{M}_\ell^\zeta$ such that an optimal solution to problem (11) is also an optimal solution to Problem (1), and vice-versa.

Coherence and Restricted Eigenvalues of a Model Matrix

Given a model matrix \mathbf{X} , [46] introduced the cumulative coherence function

$$\mu[k] := \max_{|I|=k} \max_{j \notin I} \sum_{i \in I} |\langle \mathbf{X}_j, \mathbf{X}_i \rangle|,$$

where, \mathbf{X}_j , $j = 1, \dots, p$ represent the columns of \mathbf{X} , i.e., features.

For $k = 1$, we obtain the notion of coherence introduced in [20, 19] as a measure of the maximal pairwise correlation in absolute value of the columns of \mathbf{X} :

$$\mu := \mu[1] = \max_{i \neq j} |\langle \mathbf{X}_i, \mathbf{X}_j \rangle|.$$

[14, 12] (see also [10] and references therein) introduced the notion that a matrix \mathbf{X} satisfies a restricted eigenvalue condition if

$$\lambda_{\min}(\mathbf{X}'_I \mathbf{X}_I) \geq \eta_k \quad \text{for every } I \subset \{1, \dots, p\} : |I| \leq k, \quad (12)$$

where $\lambda_{\min}(\mathbf{X}'_I \mathbf{X}_I)$ denotes the smallest eigenvalue of the matrix $\mathbf{X}'_I \mathbf{X}_I$. An inequality linking $\mu[k]$ and η_k is as follows.

Proposition 1. *The following bounds hold:*

(a) [46]: $\mu[k] \leq \mu \cdot k$.

(b) [19]: $\eta_k \geq 1 - \mu[k - 1] \geq 1 - \mu \cdot (k - 1)$.

The computations of $\mu[k]$ and η_k for general k are difficult, while μ is simple to compute. Proposition 1 provides bounds for $\mu[k]$ and η_k in terms of the coherence μ .

Operator Norms of Submatrices

The (p, q) operator norm of matrix \mathbf{A} is

$$\|\mathbf{A}\|_{p,q} := \max_{\|\mathbf{u}\|_q=1} \|\mathbf{A}\mathbf{u}\|_p.$$

We will use extensively here the $(1, 1)$ operator norm. We assume that each column vector of \mathbf{X} has unit ℓ_2 -norm. The results derived in the next proposition borrow and enhance techniques developed by [46] in the context of analyzing the ℓ_1 – ℓ_0 equivalence in compressed sensing.

Proposition 2. *For any $I \subset \{1, \dots, p\}$ with $|I| = k$ we have:*

(a) $\|\mathbf{X}'_I \mathbf{X}_I - \mathbf{I}\|_{1,1} \leq \mu[k - 1]$.

(b) *If the matrix $\mathbf{X}'_I \mathbf{X}_I$ is invertible and $\|\mathbf{X}'_I \mathbf{X}_I - \mathbf{I}\|_{1,1} < 1$, then*

$$\|(\mathbf{X}'_I \mathbf{X}_I)^{-1}\|_{1,1} \leq \frac{1}{1 - \mu[k - 1]}. \quad (13)$$

Proof

- (a) Given a set I , we define $\mathbf{G} := \mathbf{X}'_I \mathbf{X}_I - \mathbf{I}$, and let g_{ij} denote the (i, j) th entry of \mathbf{G} . For any $\mathbf{u} \in \mathbb{R}^k$ we have

$$\begin{aligned}
\max_{\|\mathbf{u}\|_1=1} \|\mathbf{G}\mathbf{u}\|_1 &= \max_{\|\mathbf{u}\|_1=1} \left(\sum_{i=1}^k \left| \sum_{j=1}^k g_{ij} u_j \right| \right) \\
&\leq \max_{\|\mathbf{u}\|_1=1} \left(\sum_{i=1}^k \sum_{j=1}^k |u_j| |g_{ij}| \right) \\
&= \max_{\|\mathbf{u}\|_1=1} \left(\sum_{j=1}^k |u_j| \sum_{i \neq j} |g_{ij}| \right) \quad (g_{jj} = 0) \\
&\leq \max_{\|\mathbf{u}\|_1=1} (\mu[k-1] \|\mathbf{u}\|_1) \quad \left(\sum_{i \neq j} |g_{ij}| \leq \mu[k-1] \right) \\
&= \mu[k-1].
\end{aligned}$$

- (b) Using $\mathbf{X}'_I \mathbf{X}_I = \mathbf{I} + \mathbf{G}$ and standard power-series convergence (which is valid since $\|\mathbf{G}\|_{1,1} < 1$) we obtain

$$\begin{aligned}
\|(\mathbf{X}'_I \mathbf{X}_I)^{-1}\|_{1,1} &= \|(\mathbf{I} + \mathbf{G})^{-1}\|_{1,1} = \sum_{i=0}^{\infty} \|\mathbf{G}\|_{1,1}^i \\
&\leq \frac{1}{1 - \|\mathbf{G}\|_{1,1}} \leq \frac{1}{1 - \mu[k-1]}. \quad \square
\end{aligned}$$

We note that Part (b) also appears in [46] for the operator norm $\|(\mathbf{X}'_I \mathbf{X}_I)^{-1}\|_{\infty, \infty}$.

Given a set $I \subset \{1, \dots, p\}$ with $|I| = k$ we let $\widehat{\boldsymbol{\beta}}_I$ denote the least squares regression coefficients obtained by regressing \mathbf{y} on \mathbf{X}_I , i.e., $\widehat{\boldsymbol{\beta}}_I = (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I \mathbf{y}$. If we append $\widehat{\boldsymbol{\beta}}_I$ with zeros in the remaining coordinates we obtain $\widehat{\boldsymbol{\beta}}$:

$$\widehat{\boldsymbol{\beta}} \in \arg \min_{\boldsymbol{\beta}: \beta_i = 0, i \notin I} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

Note that $\widehat{\boldsymbol{\beta}}$ depends on I but we will suppress the dependence on I for notational convenience.

Recall that \mathbf{X}_j , $j = 1, \dots, p$ represent the columns of \mathbf{X} ; and we will use \mathbf{x}_i , $i = 1, \dots, n$ to denote the rows of \mathbf{X} . As discussed above $\|\mathbf{X}_j\| = 1$. We order the correlations $|\langle \mathbf{X}_j, \mathbf{y} \rangle|$:

$$|\langle \mathbf{X}_{(1)}, \mathbf{y} \rangle| \geq |\langle \mathbf{X}_{(2)}, \mathbf{y} \rangle| \geq \dots \geq |\langle \mathbf{X}_{(p)}, \mathbf{y} \rangle|.$$

We finally denote by $\|\mathbf{x}_i\|_{1:k}$ the sum of the top k absolute values of the entries of x_{ij} , $j \in \{1, 2, \dots, p\}$.

Theorem 1. For any $k \geq 1$ such that $\mu[k-1] < 1$ any optimal solution $\widehat{\boldsymbol{\beta}}$ to (1) satisfies:

$$(a) \quad \|\widehat{\boldsymbol{\beta}}\|_1 \leq \frac{1}{1 - \mu[k-1]} \sum_{j=1}^k |\langle \mathbf{X}_{(j)}, \mathbf{y} \rangle|. \quad (14)$$

$$(b) \quad \|\widehat{\boldsymbol{\beta}}\|_\infty \leq \min \left\{ \frac{1}{\eta_k} \sqrt{\sum_{j=1}^k |\langle \mathbf{X}_{(j)}, \mathbf{y} \rangle|^2}, \frac{1}{\sqrt{\eta_k}} \|\mathbf{y}\|_2 \right\}. \quad (15)$$

$$(c) \quad \|\mathbf{X}\widehat{\boldsymbol{\beta}}\|_1 \leq \min \left\{ \sum_{i=1}^n \|\mathbf{x}_i\|_\infty \|\widehat{\boldsymbol{\beta}}\|_1, \sqrt{k} \|\mathbf{y}\|_2 \right\}. \quad (16)$$

$$(d) \quad \|\mathbf{X}\widehat{\boldsymbol{\beta}}\|_\infty \leq \left(\max_{i=1, \dots, n} \|\mathbf{x}_i\|_{1:k} \right) \|\widehat{\boldsymbol{\beta}}\|_\infty. \quad (17)$$

Proof

(a) Since $\widehat{\boldsymbol{\beta}}_I = (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I \mathbf{y}$ we have

$$\|\widehat{\boldsymbol{\beta}}\|_1 = \|\widehat{\boldsymbol{\beta}}_I\|_1 \leq \|(\mathbf{X}'_I \mathbf{X}_I)^{-1}\|_{1,1} \|\mathbf{X}'_I \mathbf{y}\|_1. \quad (18)$$

Note that

$$\|\mathbf{X}'_I \mathbf{y}\|_1 = \sum_{j \in I} |\langle \mathbf{X}_j, \mathbf{y} \rangle| \leq \max_{I, |I|=k} \sum_{j \in I} |\langle \mathbf{X}_j, \mathbf{y} \rangle| \leq \sum_{j=1}^k |\langle \mathbf{X}_{(j)}, \mathbf{y} \rangle|. \quad (19)$$

Applying (13) and (19) to (18), we obtain (14).

(b) We write $\widehat{\boldsymbol{\beta}}_I = \mathbf{A} \mathbf{y}$ for $\mathbf{A} = (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I$. If $\mathbf{a}_i, i = 1, \dots, k$ denote the rows of \mathbf{A} we have:

$$\|\widehat{\boldsymbol{\beta}}_I\|_\infty = \max_{i=1, \dots, k} |\langle \mathbf{a}_i, \mathbf{y} \rangle| \leq \left(\max_{i=1, \dots, k} \|\mathbf{a}_i\|_2 \right) \|\mathbf{y}\|_2. \quad (20)$$

For every $i = 1, \dots, k$ we have

$$\begin{aligned} \|\mathbf{a}_i\|_2 &\leq \max_{\|\mathbf{u}\|_2=1} \|\mathbf{A} \mathbf{u}\|_2 \\ &= \max_{\|\mathbf{u}\|_2=1} \|(\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I \mathbf{u}\|_2 \\ &\leq \lambda_{\max} \left((\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I \right) \\ &= \max \left\{ \frac{1}{d_1}, \dots, \frac{1}{d_k} \right\}, \end{aligned} \quad (21)$$

where d_1, \dots, d_k are the (nonzero) singular values of the matrix \mathbf{X}_I . To see how one arrives at (21) let us denote the singular value decomposition of $\mathbf{X}_I = \mathbf{U} \mathbf{D} \mathbf{V}'$ with $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_k)$. We then have

$$(\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I = (\mathbf{V} \mathbf{D}^{-2} \mathbf{V}') (\mathbf{U} \mathbf{D} \mathbf{V}')' = \mathbf{V} \mathbf{D}^{-1} \mathbf{U}'$$

and the singular values of $(\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I$ are thus $1/d_i$, $i = 1, \dots, k$.

The eigenvalues of $\mathbf{X}'_I \mathbf{X}_I$ are d_i^2 and from (12) we obtain that $d_i^2 \geq \eta_k$. Using (21) we thus obtain

$$\max_{i=1, \dots, k} \|\mathbf{a}_i\|_2 \leq \frac{1}{\sqrt{\eta_k}}. \quad (22)$$

Substituting the bound (22) to (20) we obtain

$$\|\widehat{\boldsymbol{\beta}}_I\|_\infty \leq \frac{1}{\sqrt{\eta_k}} \|\mathbf{y}\|_2. \quad (23)$$

Using the notation $\tilde{\mathbf{A}} = (\mathbf{X}'_I \mathbf{X}_I)^{-1}$, we have

$$\begin{aligned} \|\widehat{\boldsymbol{\beta}}_I\|_\infty &= \max_{i=1, \dots, k} |\langle \tilde{\mathbf{a}}_i, \mathbf{X}'_I \mathbf{y} \rangle| \\ &\leq \left(\max_{i=1, \dots, k} \|\tilde{\mathbf{a}}_i\|_2 \right) \|\mathbf{X}'_I \mathbf{y}\|_2 \\ &\leq \lambda_{\max} \left((\mathbf{X}'_I \mathbf{X}_I)^{-1} \right) \|\mathbf{X}'_I \mathbf{y}\|_2 \\ &= \left(\max_{i=1, \dots, k} \frac{1}{d_i^2} \right) \cdot \sqrt{\sum_{j \in I} |\langle \mathbf{X}_j, \mathbf{y} \rangle|^2} \\ &\leq \frac{1}{\eta_k} \sqrt{\sum_{j=1}^k |\langle \mathbf{X}_{(j)}, \mathbf{y} \rangle|^2}. \end{aligned} \quad (24)$$

Combining (23) and (24) we obtain (15).

(c) We have

$$\|\mathbf{X}_I \widehat{\boldsymbol{\beta}}_I\|_1 \leq \sum_{i=1}^n |\langle \mathbf{x}_i, \widehat{\boldsymbol{\beta}}_I \rangle| \leq \sum_{i=1}^n \|\mathbf{x}_i\|_\infty \|\widehat{\boldsymbol{\beta}}_I\|_1 = \sum_{i=1}^n \|\mathbf{x}_i\|_\infty \|\widehat{\boldsymbol{\beta}}_I\|_1. \quad (25)$$

Let $\mathbf{P}_I := \mathbf{X}_I (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{X}'_I$ denote the projection onto the columns of \mathbf{X}_I . We have $\|\mathbf{P}_I \mathbf{y}\|_2 \leq \|\mathbf{y}\|_2$, leading to:

$$\|\mathbf{X}_I \widehat{\boldsymbol{\beta}}_I\|_1 = \|\mathbf{P}_I \mathbf{y}\|_1 \leq \sqrt{k} \|\mathbf{P}_I \mathbf{y}\|_2 \leq \sqrt{k} \|\mathbf{y}\|_2, \quad (26)$$

where we used that for any $\mathbf{a} \in \mathbb{R}^m$, we have $\sqrt{m} \|\mathbf{a}\|_2 \geq \|\mathbf{a}\|_1$. Combining (25) and (26) we obtain (16).

(d) For any vector $\boldsymbol{\beta}_I$ which has zero entries in the coordinates outside I , we have:

$$\|\mathbf{X} \boldsymbol{\beta}_I\|_\infty \leq \max_{i=1, \dots, n} |\langle \mathbf{x}_i, \boldsymbol{\beta}_I \rangle| \leq \max_{i=1, \dots, n} \|\mathbf{x}_i\|_{1:k} \|\boldsymbol{\beta}_I\|_\infty,$$

leading to (17). □

2.4 A Simple Global Lower Bound

In this section, we show that under certain restricted eigenvalue conditions on the matrix \mathbf{X} , it is possible to compute global lower bounds to the minimum objective value of (1). Though these lower bounds can be improved by MIO techniques, they require sophisticated computational procedures—the method we present here can be computed with minimal computational effort.

Using (12) we obtain a global lower bound to (1) as follows:

$$\begin{aligned}
 \min_{\|\beta\|_0 \leq k} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 &\geq \min_{\|\beta\|_0 \leq k} (\|\mathbf{y}\|_2^2 - 2\langle \mathbf{X}'\mathbf{y}, \beta \rangle + \eta_k \|\beta\|_2^2) \\
 &= \min_{\|\beta\|_0 \leq k} \left(\eta_k \left\| \beta - \frac{1}{\eta_k} \mathbf{X}'\mathbf{y} \right\|_2^2 - \eta_k \left\| \frac{1}{\eta_k} \mathbf{X}'\mathbf{y} \right\|_2^2 + \|\beta\|_2^2 \right) \\
 &= -\frac{1}{\eta_k} \|\mathbf{H}_k(\mathbf{X}'\mathbf{y})\|_2^2 + \|\mathbf{y}\|_2^2, \tag{27}
 \end{aligned}$$

where $\mathbf{H}_k(\cdot)$ is the operator defined in (30).

Clearly (27) is a lower bound to (1) and can be computed with very little computational effort once η_k is known. Note that if $n > p$, then $\lambda_{\min}(\mathbf{X}'\mathbf{X})$ gives a lower bound to η_k (provided it is not zero). If $p > n$, $\lambda_{\min}(\mathbf{X}'\mathbf{X}) = 0$. Proposition 1(b) provides lower bounds on η_k .

3 Discrete First Order Algorithms

In this section, we develop a discrete extension of first order methods in convex optimization [40, 39] to obtain near optimal solutions for Problem (1) and its variant for the least absolute deviation (LAD) loss function. Our approach applies to the problem of minimizing any smooth convex function subject to cardinality constraints.

We will use these discrete first order methods to obtain solutions to warm start the MIO formulation. In Section 4, we will demonstrate how these methods greatly enhance the performance of the MIO.

3.1 Algorithms for Minimizing Smooth Functions Subject to Cardinality Constraints

Related Work and Contributions In the signal processing literature [8, 9] proposed iterative hard-thresholding algorithms, in the context of ℓ_0 -regularized least squares problems, i.e., Problem (4). The authors establish convergence properties of the algorithm under the assumption that \mathbf{X} satisfies coherence [8] or Restricted Isometry Property [9]. The method we propose here applies to a larger class of cardinality constrained optimization problems of the form (28), in particular, in the context of Problem (1) our algorithm and its convergence analysis do not require any form of restricted isometry property on the model matrix \mathbf{X} .

Our proposed algorithm borrows ideas from projected gradient descent methods in first order convex optimization problems [40] and generalizes it to

the discrete optimization Problem (28). We also derive new global convergence results for our proposed algorithms as presented in Theorem 2. Our proposal, with some novel modifications also applies to the non-smooth least absolute deviation loss with cardinality constraints as discussed in Section 3.3.

Consider the following optimization problem:

$$\min_{\boldsymbol{\beta}} g(\boldsymbol{\beta}) \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_0 \leq k, \quad (28)$$

where $g(\boldsymbol{\beta}) \geq 0$ is convex and has Lipschitz continuous gradient:

$$\|\nabla g(\boldsymbol{\beta}) - \nabla g(\tilde{\boldsymbol{\beta}})\| \leq \ell \|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|. \quad (29)$$

The first ingredient of our approach is the observation that when $g(\boldsymbol{\beta}) = \|\boldsymbol{\beta} - \mathbf{c}\|_2^2$ for a given \mathbf{c} , problem (28) admits a closed form solution.

Proposition 3. *An optimal solution, denoted as $\mathbf{H}_k(\mathbf{c})$, to the problem*

$$\min_{\|\boldsymbol{\beta}\|_0 \leq k} \|\boldsymbol{\beta} - \mathbf{c}\|_2^2, \quad (30)$$

can be computed as follows: $\mathbf{H}_k(\mathbf{c})$ retains the k largest (in absolute value) elements of $\mathbf{c} \in \mathbb{R}^p$ and sets the rest to zero, i.e., if $|c_{(1)}| \geq |c_{(2)}| \geq \dots \geq |c_{(p)}|$, denote the ordered values of the absolute values of the vector \mathbf{c} , then:

$$(\mathbf{H}_k(\mathbf{c}))_i = \begin{cases} c_i, & \text{if } i \in \{(1), \dots, (k)\}, \\ 0, & \text{otherwise.} \end{cases} \quad (31)$$

Proof

We provide a proof of this simple observation, for the sake of completeness.

It suffices to consider $|c_i| > 0$ for all i . Let $\boldsymbol{\beta}$ be an optimal solution to Problem (30) and let $S := \{i : \beta_i \neq 0\}$. The objective function is given by $\sum_{i \notin S} |c_i|^2 + \sum_{i \in S} (\beta_i - c_i)^2$. Note that by selecting $\beta_i = c_i$ for $i \in S$, we can make the objective function $\sum_{i \notin S} |c_i|^2$. Thus, to minimize the objective function, S must correspond to the indices of the largest k values of $|c_i|$, $i \geq 1$. \square

The operator (31) is also known as the hard-thresholding operator [18]—a notion that arises in the context of the following related optimization problem:

$$\hat{\boldsymbol{\beta}} \in \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta} - \mathbf{c}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_0, \quad (32)$$

where $\hat{\boldsymbol{\beta}}$ admits a simple closed form expression given by $\hat{\beta}_i = c_i$ if $|c_i| > \lambda$ and $\hat{\beta}_i = 0$ otherwise, for $i = 1, \dots, p$.

Remark 1. *There is an important difference between the minimizers of Problems (30) and (32). For Problem (32), the smallest (in absolute value) non-zero element in $\hat{\boldsymbol{\beta}}$ is greater than λ in absolute value. On the other hand, in Problem (30) there is no lower bound to the minimum (in absolute value) non-zero element of a minimizer, i.e., $\mathbf{H}_k(\mathbf{c})$. This needs to be taken care of using subtle techniques, while analyzing the convergence properties of Algorithm 1 (Section 3.2).*

Given a current solution β , the second ingredient of our approach is to upper bound the function $g(\eta)$ around $g(\beta)$. To do so, we use ideas from projected gradient descent methods in first order convex optimization problems [40, 39].

Proposition 4. ([40, 39]) *For a convex function $g(\beta)$ satisfying condition (29) and for any $L \geq \ell$ we have:*

$$g(\eta) \leq Q_L(\eta, \beta) := g(\beta) + \frac{L}{2} \|\eta - \beta\|_2^2 + \langle \nabla g(\beta), \eta - \beta \rangle \quad (33)$$

for all β, η with equality holding at $\beta = \eta$.

Applying Proposition 3 to the upper bound $Q_L(\eta, \beta)$ in Proposition 4 we obtain

$$\begin{aligned} & \arg \min_{\|\eta\|_0 \leq k} Q_L(\eta, \beta) \\ &= \arg \min_{\|\eta\|_0 \leq k} \left(\frac{L}{2} \left\| \eta - \left(\beta - \frac{1}{L} \nabla g(\beta) \right) \right\|_2^2 - \frac{1}{2L} \|\nabla g(\beta)\|_2^2 + g(\beta) \right) \\ &= \arg \min_{\|\eta\|_0 \leq k} \left\| \eta - \left(\beta - \frac{1}{L} \nabla g(\beta) \right) \right\|_2^2 \\ &= \mathbf{H}_k \left(\beta - \frac{1}{L} \nabla g(\beta) \right), \end{aligned} \quad (34)$$

where $\mathbf{H}_k(\cdot)$ is defined in (31). In light of (34) we are now ready to present Algorithm 1 to find a local optimal solution to problem (28).

Algorithm 1

Input: $g(\beta)$, L , ϵ .

Output: A local optimal solution β^* .

Algorithm:

1. Initialize with $\beta_1 \in \mathbb{R}^p$ such that $\|\beta_1\|_0 \leq k$.
2. For $m \geq 1$, apply (34) with $\beta = \beta_m$ to obtain β_{m+1} as:

$$\beta_{m+1} = \mathbf{H}_k \left(\beta_m - \frac{1}{L} \nabla g(\beta_m) \right) \quad (35)$$

3. Repeat Step 2, until $\|\beta_{m+1} - \beta_m\|_2 \leq \epsilon$.
4. Let $\beta_m := (\beta_{m1}, \dots, \beta_{mp})$ denote the current estimate and let $I = \text{Supp}(\beta_m) := \{i : \beta_{mi} \neq 0\}$. Solve the continuous optimization problem:

$$\min_{\beta, \beta_i=0, i \notin I} g(\beta), \quad (36)$$

and let β^* be a minimizer.

The convergence properties of Algorithm 1 are presented in Section 3.2. A variant of Algorithm 1 that has better empirical performance and uses line searches is presented next.

Algorithm 2 (with Line Search)

1. Initialize with $\beta_1 \in \mathbb{R}^p$ such that $\|\beta_1\|_0 \leq k$.
2. For $m \geq 1$,

$$\begin{aligned}\eta_m &= \mathbf{H}_k \left(\beta_m - \frac{1}{L} \nabla g(\beta_m) \right), \\ \beta_{m+1} &= \lambda_m \eta_m + (1 - \lambda_m) \beta_m,\end{aligned}\tag{37}$$

where λ_m is chosen to minimize the one-dimensional optimization problem:

$$\lambda_m \in \arg \min_{\lambda} g(\lambda \eta_m + (1 - \lambda) \beta_m).\tag{38}$$

3. Repeat Step 2, until $\|\eta_{m+1} - \eta_m\|_2 \leq \epsilon$.
4. Let η_m denote the current estimate and let $I = \text{Supp}(\eta_m)$. Solve problem (36) and let β^* be a minimizer.

Note that the iterate β_m in Algorithm 2 need not be k -sparse (i.e., need not satisfy: $\|\beta_m\|_0 \leq k$), however, η_m is k -sparse ($\|\eta_m\|_0 \leq k$). Moreover, the sequence may not lead to a decreasing set of objective values, but it satisfies:

$$g(\beta_{m+1}) \leq g(\eta_m) \not\leq g(\beta_m).$$

3.2 Convergence Analysis of Algorithm 1

In this section, we study convergence properties for Algorithm 1. Before we embark on the analysis, we need to define the notion of first order optimality for Problem (28).

Definition 1. *Given an $L \geq \ell$, the vector $\eta \in \mathbb{R}^p$ is said to be a first order stationary point of Problem (28) if $\|\eta\|_0 \leq k$ and it satisfies the following fixed point equation:*

$$\eta = \mathbf{H}_k \left(\eta - \frac{1}{L} \nabla g(\eta) \right).\tag{39}$$

We next define the notion of an ϵ -approximate first order stationary point of Problem (28):

Definition 2. *Given an $\epsilon > 0$, and $L \geq \ell$ we say that η satisfies an ϵ -approximate first order optimality condition of Problem (28) if $\|\eta\|_0 \leq k$ and*

$$\left\| \eta - \mathbf{H}_k \left(\eta - \frac{1}{L} \nabla g(\eta) \right) \right\|_2 \leq \epsilon.$$

Let $\beta_m = (\beta_{m1}, \dots, \beta_{mp})$ and $\mathbf{1}_m = (e_1, \dots, e_p)$ with $e_j = 1$, if $\beta_{mj} \neq 0$, and $e_j = 0$, if $\beta_{mj} = 0$, $j = 1, \dots, p$, i.e., $\mathbf{1}_m$ represents the sparsity pattern of the support of β_m .

Proposition 5. Consider $g(\boldsymbol{\beta})$ and ℓ as defined in (28) and (29). Let $\boldsymbol{\beta}_m, m \geq 1$ be the sequence generated by Algorithm 1. Then we have:

(a) For any $L \geq \ell$, the sequence $g(\boldsymbol{\beta}_m)$ satisfies

$$g(\boldsymbol{\beta}_m) - g(\boldsymbol{\beta}_{m+1}) \geq \frac{L - \ell}{2} \|\boldsymbol{\beta}_{m+1} - \boldsymbol{\beta}_m\|_2^2, \quad (40)$$

is decreasing and converges.

(b) If $L > \ell$, then $\boldsymbol{\beta}_{m+1} - \boldsymbol{\beta}_m \rightarrow \mathbf{0}$ as $m \rightarrow \infty$.

(c) If $L > \ell$ and $\|\liminf_{m \rightarrow \infty} \boldsymbol{\beta}_m\|_0 = k$ then the sequence $\mathbf{1}_m$ converges after finitely many iterations, i.e., there exists an iteration index M^* such that $\mathbf{1}_m = \mathbf{1}_{m+1}$ for all $m \geq M^*$. Furthermore, the sequence $\boldsymbol{\beta}_m$ is bounded and converges to a first order stationary point.

(d) If $L > \ell$ and $\|\liminf_{m \rightarrow \infty} \boldsymbol{\beta}_m\|_0 < k$, then $g(\boldsymbol{\beta}_m) \rightarrow g(\boldsymbol{\beta}^*)$ where $\boldsymbol{\beta}^* \in \arg \min g(\boldsymbol{\beta})$ is an unconstrained minimizer.

Proof

(a) Let $\boldsymbol{\beta}$ be a vector satisfying $\|\boldsymbol{\beta}\|_0 \leq k$. Using the notation $\hat{\boldsymbol{\eta}} = \mathbf{H}_k(\boldsymbol{\beta} - \frac{1}{L}\nabla g(\boldsymbol{\beta}))$ we have the following chain of inequalities:

$$\begin{aligned} g(\boldsymbol{\beta}) &= Q_L(\boldsymbol{\beta}, \boldsymbol{\beta}) \\ &\geq \inf_{\|\boldsymbol{\eta}\|_0 \leq k} Q_L(\boldsymbol{\eta}, \boldsymbol{\beta}) \\ &= \inf_{\|\boldsymbol{\eta}\|_0 \leq k} \left(\frac{L}{2} \|\boldsymbol{\eta} - \boldsymbol{\beta}\|_2^2 + \langle \nabla g(\boldsymbol{\beta}), \boldsymbol{\eta} - \boldsymbol{\beta} \rangle + g(\boldsymbol{\beta}) \right) \\ &= \inf_{\|\boldsymbol{\eta}\|_0 \leq k} \left(\frac{L}{2} \left\| \boldsymbol{\eta} - \left(\boldsymbol{\beta} - \frac{1}{L} \nabla g(\boldsymbol{\beta}) \right) \right\|_2^2 - \frac{1}{2L} \|\nabla g(\boldsymbol{\beta})\|_2^2 + g(\boldsymbol{\beta}) \right) \\ &= \left(\frac{L}{2} \|\hat{\boldsymbol{\eta}} - \left(\boldsymbol{\beta} - \frac{1}{L} \nabla g(\boldsymbol{\beta}) \right)\|_2^2 - \frac{1}{2L} \|\nabla g(\boldsymbol{\beta})\|_2^2 + g(\boldsymbol{\beta}) \right) \\ &\hspace{15em} \text{(From (34))} \\ &= \left(\frac{L}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2 + \langle \nabla g(\boldsymbol{\beta}), \hat{\boldsymbol{\eta}} - \boldsymbol{\beta} \rangle + g(\boldsymbol{\beta}) \right) \\ &= \left(\frac{L - \ell}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2 + \frac{\ell}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2 + \langle \nabla g(\boldsymbol{\beta}), \hat{\boldsymbol{\eta}} - \boldsymbol{\beta} \rangle + g(\boldsymbol{\beta}) \right) \\ &\geq \frac{L - \ell}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2 + \underbrace{\left(\frac{\ell}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2 + \langle \nabla g(\boldsymbol{\beta}), \hat{\boldsymbol{\eta}} - \boldsymbol{\beta} \rangle + g(\boldsymbol{\beta}) \right)}_{Q_\ell(\hat{\boldsymbol{\eta}}, \boldsymbol{\beta})} \\ &\geq \frac{L - \ell}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2 + g(\hat{\boldsymbol{\eta}}). \hspace{5em} \text{(From (33))} \end{aligned}$$

This chain of inequalities leads to:

$$g(\boldsymbol{\beta}) - g(\hat{\boldsymbol{\eta}}) \geq \frac{L - \ell}{2} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\beta}\|_2^2. \quad (41)$$

Applying (41) for $\boldsymbol{\beta} = \boldsymbol{\beta}_m$ and $\hat{\boldsymbol{\eta}} = \boldsymbol{\beta}_{m+1}$, the vectors generated by Algorithm 1, we obtain (40). This implies that the objective values $g(\boldsymbol{\beta}_m)$ are decreasing and since the sequence is bounded below ($g(\boldsymbol{\beta}) \geq 0$), we obtain that $g(\boldsymbol{\beta}_m)$ converges as $m \rightarrow \infty$.

(b) If $L > \ell$ and from part (a), the result follows.

(c) We begin by observing that the condition $\|\liminf_{m \rightarrow \infty} \boldsymbol{\beta}_m\|_0 = k$ is equivalent to $\liminf_{m \rightarrow \infty} \min_{i: \beta_{mi} \neq 0} |\beta_{mi}| > 0$. We next prove that the support of $\boldsymbol{\beta}_m$ converges. For the purpose of establishing of contradiction suppose that the support does not converge. Then, there are infinitely many values of m' such that $\mathbf{1}_{m'} \neq \mathbf{1}_{m'+1}$. Using the fact that $\|\boldsymbol{\beta}_m\|_0 = k$ for all large m we have

$$\|\boldsymbol{\beta}_{m'} - \boldsymbol{\beta}_{m'+1}\|_2 \geq \sqrt{\beta_{m',i}^2 + \beta_{m'+1,j}^2} \geq \frac{|\beta_{m',i}| + |\beta_{m'+1,j}|}{\sqrt{2}}, \quad (42)$$

where i, j are such that $\beta_{m'+1,i} = \beta_{m',j} = 0$. As $m' \rightarrow \infty$, the quantity in the rhs of (42) remains bounded away from zero since $\liminf_{m \rightarrow \infty} \min_{i: \beta_{mi} \neq 0} |\beta_{mi}| > 0$. This contradicts the fact that $\boldsymbol{\beta}_{m+1} - \boldsymbol{\beta}_m \rightarrow \mathbf{0}$, as established in part (b). Thus, $\mathbf{1}_m$ converges, and since $\mathbf{1}_m$ is a discrete sequence, it converges after finitely many iterations, that is $\mathbf{1}_m = \mathbf{1}_{m+1}$ for all $m \geq M^*$. Algorithm 1 becomes a vanilla gradient descent algorithm, restricted to the space $\mathbf{1}_m$ for $m \geq M^*$. Since a gradient descent algorithm for minimizing a convex function over a closed convex set leads to a sequence of iterates that converge [42, 40], we conclude that Algorithm 1 converges. Therefore, the sequence $\boldsymbol{\beta}_m$ converges to $\hat{\boldsymbol{\beta}}$, a first order stationarity point:

$$\mathbf{H}_k \left(\hat{\boldsymbol{\beta}} - \frac{1}{L} \nabla g(\hat{\boldsymbol{\beta}}) \right) = \hat{\boldsymbol{\beta}}.$$

(d) Let $\mathcal{I}_m \subset \{1, \dots, p\}$ denote the set of k largest values of the vector $(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m))$ in absolute value. By the definition of $\mathbf{H}_k(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m))$, we have

$$\left| \left(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m) \right)_i \right| \geq \left| \left(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m) \right)_j \right|,$$

for all i, j with $i \in \mathcal{I}_m$ and $j \notin \mathcal{I}_m$. Thus,

$$\begin{aligned} & \liminf_{m \rightarrow \infty} \min_{i \in \mathcal{I}_m} \left| \left(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m) \right)_i \right| \\ & \geq \liminf_{m \rightarrow \infty} \max_{j \notin \mathcal{I}_m} \left| \left(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m) \right)_j \right|. \end{aligned} \quad (43)$$

Moreover,

$$\left(\boldsymbol{\beta}_m - \mathbf{H}_k \left(\boldsymbol{\beta}_m - \frac{1}{L} \nabla g(\boldsymbol{\beta}_m) \right) \right)_i = \begin{cases} \frac{1}{L} (\nabla g(\boldsymbol{\beta}_m))_i, & i \in \mathcal{I}_m, \\ \beta_{m,i}, & \text{otherwise.} \end{cases}$$

Using the fact that $\beta_{m+1} - \beta_m \rightarrow \mathbf{0}$ we have

$$(\nabla g(\beta_m))_i \rightarrow 0, i \in \mathcal{I}_m \text{ and } \beta_{m,j} \rightarrow 0, j \notin \mathcal{I}_m$$

as $m \rightarrow \infty$. Combining with (43) we have that:

$$\begin{aligned} \liminf_{m \rightarrow \infty} \min_{i \in \mathcal{I}_m} |\beta_{mi}| &\geq \liminf_{m \rightarrow \infty} \max_{j \notin \mathcal{I}_m} \frac{1}{L} |(\nabla g(\beta_m))_j| \\ &= \frac{1}{L} \liminf_{m \rightarrow \infty} \|\nabla g(\beta_m)\|_\infty. \end{aligned}$$

Since, $\liminf_{m \rightarrow \infty} \min_{i \in \mathcal{I}_m} |\beta_{mi}| = 0$ by hypothesis, the lhs of the above inequality equals zero, which leads to $\liminf_{m \rightarrow \infty} \|\nabla g(\beta_m)\|_\infty = 0$. Thus, there is a subsequence $m' \subset \{1, 2, \dots\}$ such that $\nabla g(\beta_{m'}) \rightarrow \mathbf{0}$, i.e., $\liminf_{m \rightarrow \infty} \nabla g(\beta_m) \rightarrow \mathbf{0}$. Since, $g(\beta_m)$ is a decreasing sequence, this implies that $g(\beta_m) \rightarrow g(\beta^*)$, where, β^* is an unconstrained (without cardinality constraints) solution to $\min g(\beta)$. \square

Proposition 5 establishes that Algorithm 1 either converges to a first order stationarity point (part (c)) or it converges to a global optimal solution (part (d)), but does not quantify the rate of convergence. We next characterize the rate of convergence of the algorithm to an ϵ -approximate first order stationary point.

Theorem 2. *Let $L > \ell$ and $\hat{\beta}$ denote a first order stationary point of Algorithm 1. After M iterations Algorithm 1 satisfies*

$$\min_{m=1, \dots, M} \|\beta_{m+1} - \beta_m\|_2^2 \leq \frac{2(g(\beta_1) - g(\hat{\beta}))}{M(L - \ell)}, \quad (44)$$

where $g(\beta_m) \downarrow g(\hat{\beta})$ as $m \rightarrow \infty$.

Proof

Summing inequalities (40) for $1 \leq m \leq M$. we obtain

$$\sum_{m=1}^M (g(\beta_m) - g(\beta_{m+1})) \geq \frac{L - \ell}{2} \sum_{m=1}^M \|\beta_{m+1} - \beta_m\|_2^2, \quad (45)$$

leading to

$$g(\beta_1) - g(\beta_{M+1}) \geq \frac{M(L - \ell)}{2} \min_{m=1, \dots, M} \|\beta_{m+1} - \beta_m\|_2^2.$$

Since the decreasing sequence $g(\beta_{m+1})$ converges to $g(\hat{\beta})$ by Proposition 5 we have:

$$\frac{g(\beta_1) - g(\hat{\beta})}{M} \geq \frac{g(\beta_1) - g(\beta_{M+1})}{M} \geq \frac{(L - \ell)}{2} \min_{m=1, \dots, M} \|\beta_{m+1} - \beta_m\|_2^2.$$

Theorem 2 implies that for any $\epsilon > 0$ there exists $M = O(\frac{1}{\epsilon})$ such that for some $1 \leq m^* \leq M$

$$\|\beta_{m^*+1} - \beta_{m^*}\|_2^2 \leq \epsilon.$$

Polishing Coefficients on the Active Set

Algorithm 1 *detects* the active set after a few iterations. Once the active set stabilizes, the algorithm may take a number of iterations to estimate the values of the regression coefficients on the active set to a high accuracy level.

In this context, we found the following simple polishing of coefficients to be useful. When the algorithm has converged to a tolerance of ϵ ($\approx 10^{-4}$), we fix the current active set, \mathcal{I} , and solve the following lower-dimensional convex optimization problem:

$$\min_{\beta, \beta_i=0, i \notin \mathcal{I}} g(\beta). \quad (46)$$

In the context of the least squares and the least absolute deviation problems, the optimization Problem (46) reduces to a smaller dimensional least squares and a linear optimization problem respectively, which can be solved very efficiently up to a very high level of accuracy.

3.3 Application to Least Squares

For the support constrained problem with squared error loss, we have

$$g(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2, \quad \nabla g(\beta) = -\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta)$$

The general algorithmic framework developed above applies in a straightforward fashion for this special case. Note that for this case $\ell = \lambda_{\max}(\mathbf{X}'\mathbf{X})$.

The polishing of the regression coefficients in the active set can be performed via a least squares problem on \mathbf{y} , \mathbf{X}_I , where I denotes the support of the regression coefficients.

3.4 Application to Least Absolute Deviation

We will now show how the method proposed in the previous section applies to the least absolute deviation problem with support constraints in β :

$$\begin{aligned} \min_{\beta} \quad & g_1(\beta) := \|\mathbf{Y} - \mathbf{X}\beta\|_1 \\ \text{s.t.} \quad & \|\beta\|_0 \leq k. \end{aligned} \quad (47)$$

Since $g_1(\beta)$ is non-smooth, our framework does not apply directly. We smooth the non-differentiable $g_1(\beta)$ so that we can apply Algorithms 1 and 2. Following [38] we make use of the following min-max representation of $g_1(\beta)$:

$$\begin{aligned} g_1(\beta) = \sup_{\mathbf{w}} \quad & \langle \mathbf{Y} - \mathbf{X}\beta, \mathbf{w} \rangle \\ \text{s.t.} \quad & \|\mathbf{w}\|_{\infty} \leq 1 \end{aligned} \quad (48)$$

and perturb the linear functional in (48) as follows:

$$\begin{aligned} g_1(\beta; \tau) = \sup_{\mathbf{w}} \quad & \langle \mathbf{Y} - \mathbf{X}\beta, \mathbf{w} \rangle - \frac{\tau}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & \|\mathbf{w}\|_{\infty} \leq 1. \end{aligned} \quad (49)$$

The following properties follow from [38]:

$$-\frac{\tau n}{2} \leq g_1(\beta; \tau) - g_1(\beta) \leq \frac{\tau n}{2}, \quad \forall \beta. \quad (50)$$

$$\|\nabla g_1(\boldsymbol{\beta}; \tau) - \nabla g_1(\tilde{\boldsymbol{\beta}}; \tau)\| \leq \frac{\lambda_{\max}(\mathbf{X}'\mathbf{X})}{\tau} \|\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}\|. \quad (51)$$

The parameter τ controls the tightness of approximation (50) and also the smoothness (51) of $g_1(\boldsymbol{\beta}; \tau)$. In order to obtain a good approximation to Problem (47), we found the following strategy to be useful in practice:

1. Fix $\tau > 0$, initialize with $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ and repeat the following steps [2]–[3] till convergence:
2. Apply Algorithm 1 (or Algorithm 2) to the smooth function $g_1(\boldsymbol{\beta}; \tau)$. Let $\boldsymbol{\beta}_\tau^*$ be the limiting solution.
3. Decrease $\tau \leftarrow \tau\gamma$ for some pre-defined constant $\gamma = 0.8$ (say), and go back to step [1] with $\boldsymbol{\beta}_0 = \boldsymbol{\beta}_\tau^*$. Exit if $\tau < \text{TOL}$, for some pre-defined tolerance.

4 Computational Experiments for Subset Selection with Least Squares Loss

In this section, we present a variety of computational tests to assess the algorithmic and statistical performances of our approach. We consider both the classical overdetermined case with $n > p$ (Section 4.2) and the high dimensional $p \gg n$ case (Section 4.3) for the least squares loss function with support constraints.

4.1 Description of Experimental Data

We demonstrate the performance of our proposal via a series of experiments on both synthetic and real data.

Synthetic Datasets. We consider a collection of problems where $\mathbf{x}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma})$, $i = 1, \dots, n$ are independent realizations from a p -dimensional multivariate normal distribution with mean zero and covariance matrix $\boldsymbol{\Sigma} := (\sigma_{ij})$. The columns of the \mathbf{X} matrix were subsequently standardized to have unit ℓ_2 norm. For a fixed $\mathbf{X}_{n \times p}$, we generated the response \mathbf{y} as follows: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^0 + \boldsymbol{\epsilon}$, where $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. We denote the number of nonzeros in $\boldsymbol{\beta}^0$ by k_0 . The choice of $\mathbf{X}, \boldsymbol{\beta}^0, \sigma$ determines the Signal-to-Noise Ratio (SNR) of the problem, which is defined as:

$$\text{SNR} = \frac{\text{var}(\mathbf{x}'\boldsymbol{\beta}^0)}{\sigma^2}.$$

We considered the following four different examples:

Example 1: We took $\sigma_{ij} = \rho^{|i-j|}$ for $i, j \in \{1, \dots, p\} \times \{1, \dots, p\}$. We consider different values of $k_0 \in \{5, 10\}$ and $\beta_i^0 = 1$ for k_0 equi-spaced values² of i in the range $\{1, 2, \dots, p\}$.

²In the case where exactly equi-spaced values are not possible we rounded the indices to the nearest large integer value.

Example 2: We took $\Sigma = \mathbf{I}_{p \times p}$, $k_0 = 5$ and $\beta^0 = (\mathbf{1}'_{5 \times 1}, \mathbf{0}'_{p-5 \times 1})' \in \mathbb{R}^p$.

Example 3: We took $\Sigma = \mathbf{I}_{p \times p}$, $k_0 = 10$ and $\beta_i^0 = \frac{1}{2} + (10 - \frac{1}{2}) \frac{(i-1)}{k_0}$, $i = 1, \dots, 10$ and $\beta_i^0 = 0, \forall i > 10$ — i.e., a vector with ten nonzero entries, with the nonzero values being equally spaced in the interval $[\frac{1}{2}, 10]$.

Example 4: We took $\Sigma = \mathbf{I}_{p \times p}$, $k_0 = 6$ and $\beta^0 = (-10, -6, -2, 2, 6, 10, \mathbf{0}_{p-6})$, i.e., a vector with six nonzero entries, equally spaced in the interval $[-10, 10]$.

Real Datasets. We considered the Diabetes dataset analyzed in [21]. We used the dataset with all the second order interactions included in the model, which resulted in 64 predictors. We reduced the sample size to $n = 350$ by taking a random sample and standardized the response and the columns of the model matrix to have zero means and unit ℓ_2 -norm.

In addition to the above, we also considered a real microarray dataset, the Leukemia data [16]. We downloaded the processed dataset from <http://stat.ethz.ch/~dettling/bagboost.html>, which had $n = 72$ binary responses and more than 3000 predictors. We standardized the response and columns of features to have zero means and unit ℓ_2 -norm. We reduced the set features to 1000 by retaining the features maximally correlated (in absolute value) to the response. We call the resulting feature matrix $\mathbf{X}_{n \times p}$ with $n = 72, p = 1000$. We then generated a semi-synthetic dataset with continuous response as $\mathbf{y} = \mathbf{X}\beta^0 + \epsilon$, where the first five coefficients of β^0 were taken as one and the rest as zero. The noise was distributed as $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$, with σ^2 chosen to get a SNR=7.

Computer Specifications and Software Computations were carried out in a linux 64 bit server—Intel(R) Xeon(R) eight-core processor @ 1.80GHz, 16 GB of RAM for the overdetermined $n > p$ case and in a Dell Precision T7600 computer with an Intel Xeon E52687 sixteen-core processor @ 3.1GHz, 128GB of Ram for the high-dimensional $p \gg n$ case. The discrete first order methods were implemented in MATLAB 2012b. We used Gurobi [29] version 5.5 and the MATLAB interface to Gurobi for all of our experiments, apart from the computations for synthetic data for $n > p$, which were done in Gurobi via its Python 2.7 interface.

4.2 The Overdetermined Regime: $n > p$

Using the Diabetes dataset and synthetic datasets, we demonstrate the combined effect of using the discrete first order methods with the MIO approach. Together, these methods show improvements in obtaining good upper bounds and in closing the MIO gap to certify global optimality. Using synthetic datasets where we know the true linear regression model, we perform side-by-side comparisons of this method with several other state-of-the-art algorithms designed to estimate sparse linear models.

4.2.1 Obtaining Good Upper Bounds

We conducted experiments to evaluate the performance of our methods in terms of obtaining high quality solutions for Problem (1).

We considered the following three algorithms:

- (a) Algorithm 2 with fifty random initializations³. We took the solution corresponding to the best objective value.
- (b) MIO with cold start, i.e., formulation (9) with a time limit of 500 seconds.
- (c) MIO with warm start. This was the MIO formulation initialized with the discrete first order optimization solution obtained from (a). This was run for a total of 500 seconds.

To compare the different algorithms in terms of the quality of upper bounds, we run for every instance all the algorithms and obtain the best solution among them, say, f_* . If f_{alg} denotes the value of the best subset objective function for method “alg”, then we define the relative accuracy of the solution obtained by “alg” as:

$$\text{Relative Accuracy} = (f_{\text{alg}} - f_*)/f_*, \quad (52)$$

where $\text{alg} \in \{(a), (b), (c)\}$ as described above.

We did experiments for the Diabetes dataset for different values of k (see Table 1). For each of the algorithms we report the amount of time taken by the algorithm to reach the best objective value during the time of 500 seconds.

k	Discrete First Order		MIO Cold Start		MIO Warm Start	
	Accuracy	Time	Accuracy	Time	Accuracy	Time
9	0.1306	1	0.0036	500	0	346
20	0.1541	1	0.0042	500	0	77
49	0.1915	1	0.0015	500	0	87
57	0.1933	1	0	500	0	2

Table 1: Quality of upper bounds for Problem (1) for the Diabetes dataset, for different values of k . We see that the MIO equipped with warm starts deliver the best upper bounds in the shortest overall times. The run time for the MIO with warm start includes the time taken by the discrete first order method (which were all less than a second).

Using the discrete first order methods in combination with the MIO algorithm resulted in finding the best possible relative accuracy in a matter of a few minutes.

³we took fifty random starting values around $\mathbf{0}$ of the form $\min(i - 1, 1)\epsilon$, $i = 1, \dots, 50$, where $\epsilon \sim N(\mathbf{0}_{p \times 1}, \mathbf{4I})$. We found empirically that Algorithm 2 provided better upper bounds than Algorithm 1.

4.2.2 Improving MIO Performance via Warm Starts

We performed a series of experiments on the Diabetes dataset to obtain a globally optimal solution to Problem (1) via our approach and to understand the implications of using advanced warm starts to the MIO formulation in terms of certifying optimality. For each choice of k we ran Algorithm 2 with fifty random initializations. They took less than a few seconds to run. We used the best solution as an advanced warm start to the MIO formulation (9). For each of these examples, we also ran the MIO formulation without any warm start (we refer to this as “cold start”). Figure 3 summarizes the results. The figure shows that in the presence of warm starts, the MIO closes the optimality gap significantly faster than those without advanced warm starts.

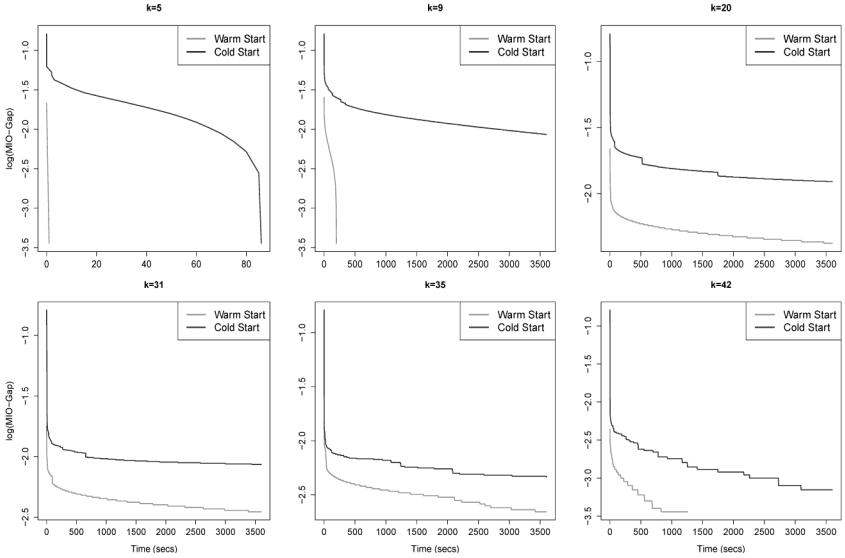


Figure 3: The evolution of the MIO optimality gap (in $\log_{10}(\cdot)$ scale) for Problem (1), for the Diabetes dataset with $n = 350, p = 64$ with and without warm starts for different values of k . The MIO significantly benefits by advanced warm starts delivered by Algorithm 2. In all of these examples, the global optimum was found within a very small fraction of the total time, but the proof of global optimality came later. As the number of possible solutions grows as $\binom{p}{k}$, it takes longer to prove optimality for $k = 31, 35$ compared to $k = 42$.

4.2.3 Statistical Performance

We considered datasets as described in Example 1, Section 4.1—we took different values of n, p with $n > p$, ρ with $k_0 = 10$.

Competing Methods and Performance Measures For every example, we considered the following learning procedures for comparison purposes: (a) the MIO approach equipped warm starts from Algorithm 2 (annotated as

“MIO” in the figure), (b) the Lasso, (c) Sparsenet and (d) stepwise regression (annotated as “Step” in the figure).

We used R to compute Lasso, Sparsenet and stepwise regression using the glmnet 1.7.3, Sparsenet and Stats 3.0.2 packages respectively, which were all downloaded from CRAN at <http://cran.us.r-project.org/>.

For each procedure, we obtained the “optimal” tuning parameter by selecting the model that achieved the best predictive performance on a held out validation set. Once the model $\hat{\beta}$ was selected, we obtained the prediction error as:

$$\text{Prediction Error} = \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta^0\|_2^2 / \|\mathbf{X}\beta^0\|_2^2. \quad (53)$$

We report “prediction error” and number of non-zeros in the optimal model in our results. The results were averaged over ten random instances, for different realizations of \mathbf{X}, ϵ . For every run: the training and validation data had a fixed \mathbf{X} but random noise ϵ .

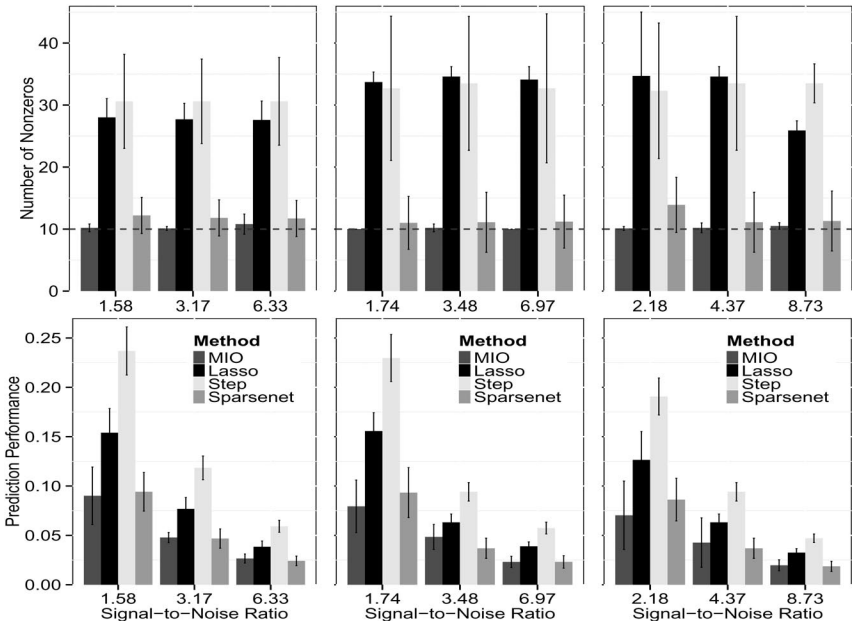


Figure 4: Figure showing the sparsity (upper panel) and predictive performances (bottom panel) for different subset selection procedures for the least squares loss. Here, we consider data generated as per Example 1, with $n = 500, p = 100, k_0 = 10$, for three different SNR values with [Left Panel] $\rho = 0.5$, [Middle Panel] $\rho = 0.8$, and [Right Panel] $\rho = 0.9$. The dashed line in the top panel represents the true number of nonzero values. For each of the procedures, the optimal model was selected as the one which produced the best prediction accuracy on a separate validation set, as described in Section 4.2.3.

Figure 4 presents results for data generated as per Example 1 with $n = 500$ and $p = 100$. We see that the MIO procedure performs very well across all the

examples. Among the methods, MIO performs the best, followed by Sparsenet, Lasso with Step(wise) exhibiting the worst performance. In terms of prediction error, the MIO performs the best, only to be marginally outperformed by Sparsenet in a few instances. This further illustrates the importance of using non-convex methods in sparse learning. Note that the MIO approach, unlike Sparsenet certifies global optimality in terms of solving Problem 1. However, based on the plots in the upper panel, Sparsenet selects a few redundant variables unlike MIO. Lasso delivers quite dense models and pays the price in predictive performance too, by selecting wrong variables. As the value of SNR increases, the predictive power of the methods improve, as expected. The differences in predictive errors between the methods diminish with increasing SNR values. With increasing values of ρ (from left panel to right panel in the figure), the number of non-zeros selected by the Lasso in the optimal model increases.

We also performed experiments with $n = 1000, p = 50$ for data generated as per Example 1. We solved the problems to provable optimality and found that the MIO performed very well when compared to other competing methods. We do not report the experiments for brevity.

4.2.4 MIO Model Training

We trained a sequence of best subset models (indexed by k) by applying the MIO approach with warm starts. Instead of running the MIO solvers from scratch for different values of k , we used *callbacks*, a feature of integer optimization solvers. Callbacks allow the user to solve an initial model, and then add additional constraints to the model one at a time. These “cuts” reduce the size of the feasible region without having to rebuild the entire optimization model. Thus, in our case, we can save time by building the initial optimization model for $k = p$. Once the solution for $k = p$ is obtained, a cut can be added to the model: $\sum_{i=1}^p z_i \leq k$ for $k = p - 1$ and the model can be re-solved from this point. We apply this procedure until we arrive at a model with $k = 1$.

For each value of k tested, the MIO best subset algorithm was set to stop the first time either an optimality gap of 1% was reached or a time limit of 15 minutes was reached. Additionally, we only tested values of k from 5 through 25, and used Algorithm 2 to warm start the MIO algorithm. We observed that it was possible to obtain speedups of a factor of 2-4 by carefully tuning the optimization solver for a particular problem, but chose to maintain generality by solving with default parameters. Thus, we do not report times with the intention of accurately benchmarking the best possible time but rather to show that it is computationally tractable to solve problems to optimality using modern MIO solvers.

4.3 The High-Dimensional Regime: $p \gg n$

In this section, we investigate

- (a) the evolution of upper bounds in the high-dimensional regime,

- (b) the effect of a bounding box formulation on the speed of closing the optimality gap,
- (c) the statistical performance of the MIO approach in comparison to other state-of-the art methods

4.3.1 Obtaining Good Upper Bounds

We performed tests similar to those in Section 4.2.1 for the $p \gg n$ regime. We tested a synthetic dataset corresponding to Example 2 with $n = 30, p = 2000$ for varying SNR values (see Table 2) over a time of 500s. As before, using the discrete first order methods in combination with the MIO algorithm resulted in finding the best possible upper bounds in the shortest possible times.

	k	Discrete First Order		MIO Cold Start		MIO Warm Start	
		Accuracy	Time	Accuracy	Time	Accuracy	Time
SNR = 3	5	0.1647	37.2	1.0510	500	0	72.2
	6	0.6152	41.1	0.2769	500	0	77.1
	7	0.7843	40.7	0.8715	500	0	160.7
	8	0.5515	38.8	2.1797	500	0	295.8
	9	0.7131	45.0	0.4204	500	0	96.0
SNR = 7	5	0.5072	45.6	0.7737	500	0	65.6
	6	1.3221	40.3	0.5121	500	0	82.3
	7	0.9745	40.9	0.7578	500	0	210.9
	8	0.8293	40.5	1.8972	500	0	262.5
	9	1.1879	44.2	0.4515	500	0	254.2

Table 2: The quality of upper bounds for Problem (1) obtained by Algorithm 2, MIO with cold start and MIO warm-started with Algorithm 2. We consider the synthetic dataset of Example 2 with $n = 30, p = 2000$ and different values of SNR. The MIO method, when warm-started with the first order solution performs the best in terms of getting a good upper bound in the shortest time. The metric “Accuracy” is defined in (52). The first order methods are fast but need not lead to highest quality solutions on their own. MIO improves the quality of upper bounds delivered by the first order methods and their combined effect leads to the best performance.

We also did experiments on the Leukemia dataset. In Figure 5 we demonstrate the evolution of the objective value of the best subset problem for different values of k . For each value of k , we warm-started the MIO with the solution obtained by Algorithm 2 and allowed the MIO solver to run for 4000 seconds. The best objective value obtained at the end of 4000 seconds is denoted by f_* . We plot the Relative Accuracy, i.e., $(f_t - f_*)/f_*$, where f_t is the objective value obtained after t seconds. The figure shows that the solution obtained by Algorithm 2 is improved by the MIO on various instances and the time taken to improve the upper bounds depends upon k . In general, for smaller values of k the upper bounds obtained by the MIO algorithm stabilize earlier, i.e., the MIO finds improved solutions faster than for larger values of k .

Leukemia data-set

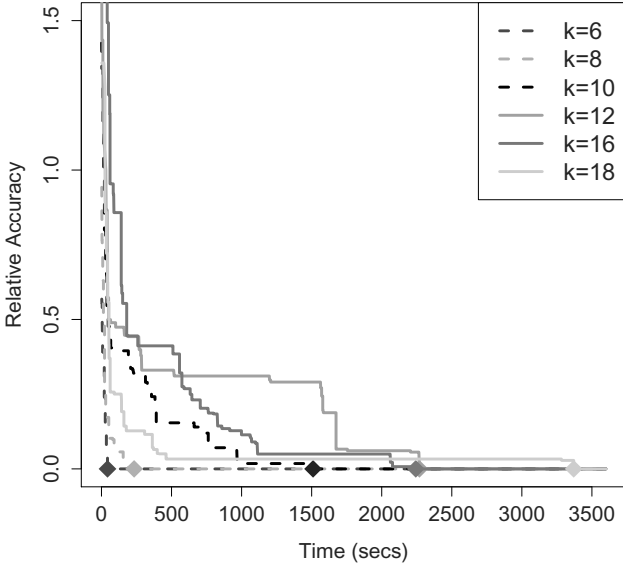


Figure 5: Behavior of MIO aided with warm start in obtaining good upper bounds over time for the Leukemia dataset ($n = 72, p = 1000$). The vertical axis shows relative accuracy, i.e., $(f_t - f_*)/f_*$, where f_t is the objective value obtained after t seconds and f_* denotes the best objective value obtained by the method after 4000 seconds. The colored diamonds correspond to the locations where the MIO (with warm start) attains the best solution. The figure shows that MIO improves the solution obtained by the first order method in all the instances. The time at which the best possible upper bound is obtained depends upon the choice of k . Typically larger k values make the problem harder—hence the best solutions are obtained after a longer wait.

4.3.2 Bounding Box Formulation

With the aid of advanced warm starts as provided by Algorithm 2, the MIO obtains a very high quality solution very quickly—in most of the examples the solution thus obtained turns out to be the global minimum. However, in the typical “high-dimensional” regime, with $p \gg n$, we observe that the certificate of global optimality comes later as the lower bounds of the problem “evolve” slowly. This is observed even in the presence of warm starts and using the implied bounds as developed in Section 2.2 and is aggravated for the cold-started MIO formulation (10).

To address this, we consider the MIO formulation (54) obtained by adding bounding boxes around a local solution. These restrictions *guide* the MIO in restricting its *search* space and enable the MIO to certify global optimality inside that bounding box. We consider the following additional bounding box constraints to the MIO formulation (10):

$$\left\{ \beta : \|\mathbf{X}\beta - \mathbf{X}\beta_0\|_1 \leq \mathcal{L}_{\ell, \text{loc}}^\zeta \right\} \cap \left\{ \beta : \|\beta - \beta_0\|_1 \leq \mathcal{L}_{\ell, \text{loc}}^\beta \right\},$$

where, β_0 is a candidate sparse solution. The radii of the two ℓ_1 -balls above, namely, $\mathcal{L}_{\ell,\text{loc}}^\zeta$ and $\mathcal{L}_{\ell,\text{loc}}^\beta$ are user-defined parameters and control the size of the feasible set.

Using the notation $\zeta = \mathbf{X}\beta$ we have the following MIO formulation (equipped with the additional bounding boxes):

$$\begin{aligned}
& \min_{\beta, \mathbf{z}, \zeta} && \frac{1}{2} \zeta^T \zeta - \langle \mathbf{X}'\mathbf{y}, \beta \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 \\
& \text{s.t.} && \zeta = \mathbf{X}\beta \\
& && (\beta_i, 1 - z_i) : \text{SOS type-1}, \quad i = 1, \dots, p \\
& && z_i \in \{0, 1\}, \quad i = 1, \dots, p \\
& && \sum_{i=1}^p z_i \leq k \\
& && -\mathcal{M}_U \leq \beta_i \leq \mathcal{M}_U, \quad i = 1, \dots, p \\
& && \|\beta\|_1 \leq \mathcal{M}_\ell \\
& && -\mathcal{M}_U^\zeta \leq \zeta_i \leq \mathcal{M}_U^\zeta, \quad i = 1, \dots, n \\
& && \|\zeta\|_1 \leq \mathcal{M}_\ell^\zeta \\
& && \|\zeta - \zeta_0\|_1 \leq \mathcal{L}_{\ell,\text{loc}}^\zeta \\
& && \|\beta - \beta_0\|_1 \leq \mathcal{L}_{\ell,\text{loc}}^\beta.
\end{aligned} \tag{54}$$

For large values of $\mathcal{L}_{\ell,\text{loc}}^\zeta$ (respectively, $\mathcal{L}_{\ell,\text{loc}}^\beta$) the constraints on $\mathbf{X}\beta$ (respectively, β) become ineffective and one gets back formulation (10). To see the impact of these additional cutting planes in the MIO formulation, we consider a few examples as illustrated in Figures 6,7,8.

Interpretation of the Bounding Boxes. A local bounding box in the variable $\zeta = \mathbf{X}\beta$ directs the MIO solver to seek for candidate solutions that deliver models with predictive accuracy “similar” (controlled by the radius of the ball) to a reference predictive model, given by ζ_0 . In our experiments, we typically chose ζ_0 as the solution delivered by running MIO (warm-started with a first order solution) for a few hundred to a few thousand seconds. More generally, ζ_0 may be selected by any other sparse learning method. In our experiments, we found that the run-time behavior of the MIO depends upon how correlated the columns of \mathbf{X} are.

Similarly, a bounding box around β directs the MIO to look for solutions in the neighborhood of a reference point β_0 . In our experiments, we chose the reference β_0 as the solution obtained by MIO (warm-started with a first order solution) and allowing it to run for a few hundred to a few thousand seconds. We observed in our experiments that the MIO solver in presence of bounding boxes in the β -space certified optimality and in the process finding better solutions; much faster than the ζ -bounding box method.

Note that the β -bounding box constraint leads to $O(p)$ and the ζ -box leads to $O(n)$ constraints. Thus, when $p \gg n$ the additional ζ constraints add a fewer number of extra variables when compared to the β constraints.

Experiments In the first set of experiments, we consider the Leukemia dataset with $n = 72$, $p = 1000$. We took two different values of $k \in \{5, 10\}$ and for each case we ran Algorithm 2 with several random restarts. The best solution thus obtained was used to warm start the MIO formulation (10), which we ran for an additional 3600 seconds. The solution thus obtained is denoted by β_0 . We then consider formulation (54) with $\mathcal{L}_{\ell, \text{loc}}^\zeta = \infty$ and different values of $\mathcal{L}_{\ell, \text{loc}}^\beta = \text{Frac}$ (as annotated in Figure 6) — the results are displayed in Figure 6.

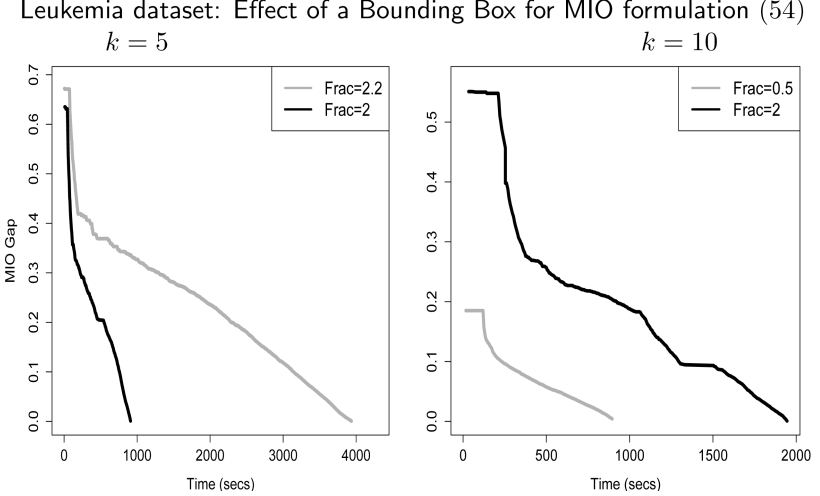


Figure 6: The effect of the MIO formulation (54) for the Leukemia dataset, for different values of k . Here $\mathcal{L}_{\ell, \text{loc}}^\zeta = \infty$ and $\mathcal{L}_{\ell, \text{loc}}^\beta = \text{Frac}$. For each value of k , the global minimum obtained was the same for the different choices of $\mathcal{L}_{\ell, \text{loc}}^\beta$.

We consider another set of experiments to demonstrate the performance of the MIO in certifying global optimality for different synthetic datasets with varying n, p, k as well as with different structures on the bounding box. In the first case, we generated data as per Example 1 with $\rho = 0.9$, $k_0 = 5$. We consider the case with $\zeta_0 = \mathbf{X}\beta_0$, $\mathcal{L}_{\ell, \text{loc}}^\beta = \infty$ and $\mathcal{L}_{\ell, \text{loc}}^\zeta = 0.5\|\mathbf{X}\beta_0\|_1$, where β_0 is a k -sparse solution obtained from the MIO formulation (10) run with a time limit of 1000 seconds, after being warm-started with Algorithm 2. The results are displayed in Figure 7 [Left Panel]. In the second case (with data same as before) we obtained β_0 in the same fashion as described before—we took a bounding box around β_0 , and left the box constraint around $\mathbf{X}\beta_0$ inactive, i.e., we set $\mathcal{L}_{\ell, \text{loc}}^\zeta = \infty$ and $\mathcal{L}_{\ell, \text{loc}}^\beta = \|\beta_0\|_1/k$. We performed two sets of experiments, where the data were generated based on different SNR values—the results are displayed in Figure 7 with SNR=1 [Middle Panel] and SNR = 3 [Right Panel].

In the same vein, we have Figure 8 studying the effect of formulations (54) for synthetic datasets generated as per Example 1 with $n = 50$, $p = 1000$, $\rho = 0.9$ and $k_0 = 5$.

Evolution of the MIO gap for (54), effect of type of bounding box
($n = 50, p = 500$).

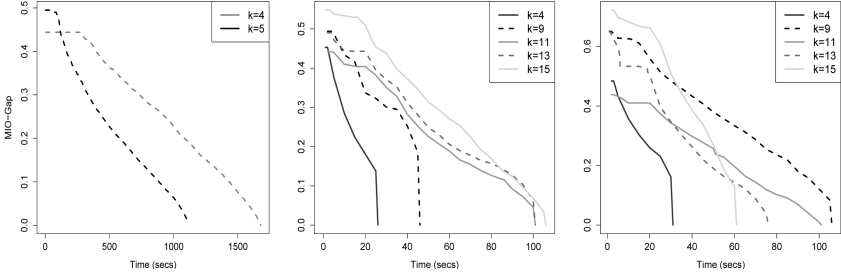


Figure 7: The effect of the MIO formulation (54) for a synthetic dataset as in Example 1 with $\rho = 0.9, k_0 = 5, n = 50, p = 500$, for different values of k . [Left Panel] $\mathcal{L}_{\ell, \text{loc}}^{\zeta} = 0.5 \|\mathbf{X}\beta_0\|_1$ and $\mathcal{L}_{\ell, \text{loc}}^{\beta} = \infty$ for a data-set with SNR = 3. [Middle Panel] $\mathcal{L}_{\ell, \text{loc}}^{\zeta} = \infty, \mathcal{L}_{\ell, \text{loc}}^{\beta} = \|\beta_0\|_1/k$ and SNR = 1. [Right Panel] $\mathcal{L}_{\ell, \text{loc}}^{\zeta} = \infty, \mathcal{L}_{\ell, \text{loc}}^{\beta} = \|\beta_0\|_1/k$ and SNR = 3. The figure shows that the bounding boxes in terms of $\mathbf{X}\beta$ (left-panel) make the problem harder to solve, when compared to bounding boxes around β (middle and right panels). A possible reason is due to the strong correlations among the columns of \mathbf{X} . The SNR values do not seem to have a big impact on the run-times of the algorithms (middle and right panels).

4.3.3 Statistical Performance

To understand the statistical behavior of MIO when compared to other approaches for learning sparse models, we considered synthetic datasets for values of n ranging from 30 – 50 and values of p ranging from 1000 – 2000. The following methods were used for comparison purposes

- (a) Algorithm 2. Here we used fifty different random initializations around $\mathbf{0}$, of the form $\min(i - 1, 1)N(\mathbf{0}_{p \times 1}, 4\mathbf{I}), i = 1, \dots, 50$ and took the solution corresponding to the best objective value.
- (b) The MIO approach with warm starts from part (a).
- (c) The Lasso solution.
- (d) The Sparsenet solution.

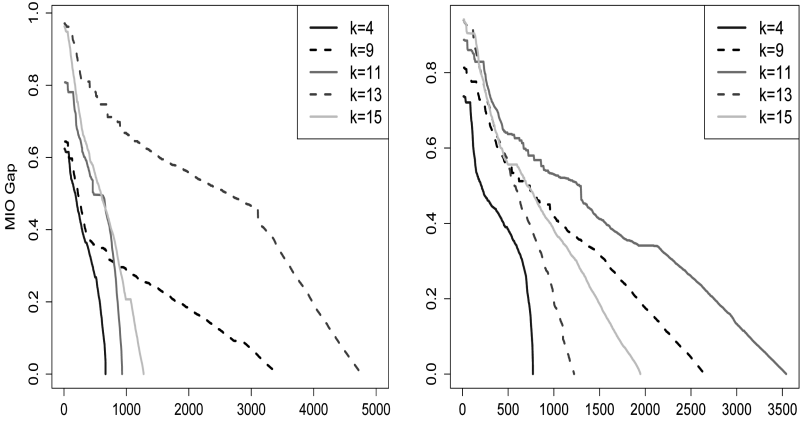
For methods (a), (b) we considered ten equi-spaced values of k in the range $[3, 2k_0]$ (including the optimal value of k_0). For each of the methods, the best model was selected in the same fashion as described in Section 4.2.3 using separate validation sets.

In Figure 9 and Figure 10 we present selected representative results from four different examples described in Section 4.1.

In Figure 9 the left panel shows the performance of different methods for Example 1 with $n = 50, p = 1000, \rho = 0.8, k_0 = 5$. In this example, there are five non-zero coefficients: the features corresponding to the non-zero coefficients are weakly correlated and a feature having a non-zero coefficient is highly

Evolution of the MIO gap for (54), effect of bounding box radii
 ($n = 50, p = 1000$).

$$\mathcal{L}_{\ell, \text{loc}}^{\zeta} = \infty \text{ and } \mathcal{L}_{\ell, \text{loc}}^{\beta} = 2\|\beta_0\|_1/k$$



$$\mathcal{L}_{\ell, \text{loc}}^{\zeta} = \infty \text{ and } \mathcal{L}_{\ell, \text{loc}}^{\beta} = \|\beta_0\|_1/k$$

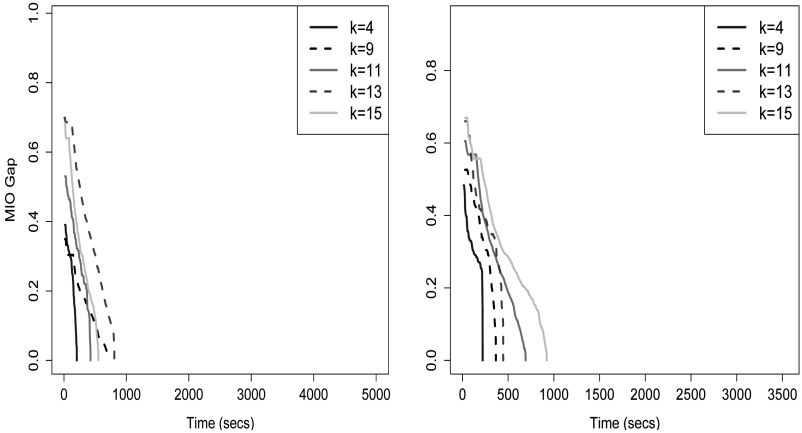


Figure 8: The evolution of the MIO gap with varying radii of bounding boxes for MIO formulation (54). The top panel has radii twice the size of the bottom panel. The dataset considered is generated as per Example 1 with $n = 50, p = 1000, \rho = 0.9$ and $k_0 = 5$ for different values of SNR: [Left Panel] SNR = 1, [Right Panel] SNR = 3. For each case, different values of k have been considered. The top panel has a bounding box radii which is twice the corresponding case in the lower panel. As expected, the times for the MIO gaps to close depends upon the radii of the boxes. The optimal solutions obtained were found to be insensitive to the choice of the bounding box radius.

correlated with a feature having a zero coefficient. In this situation, the Lasso selects a very dense model since it fails to distinguish between a zero and a non-zero coefficient when the variables are correlated—it brings both the

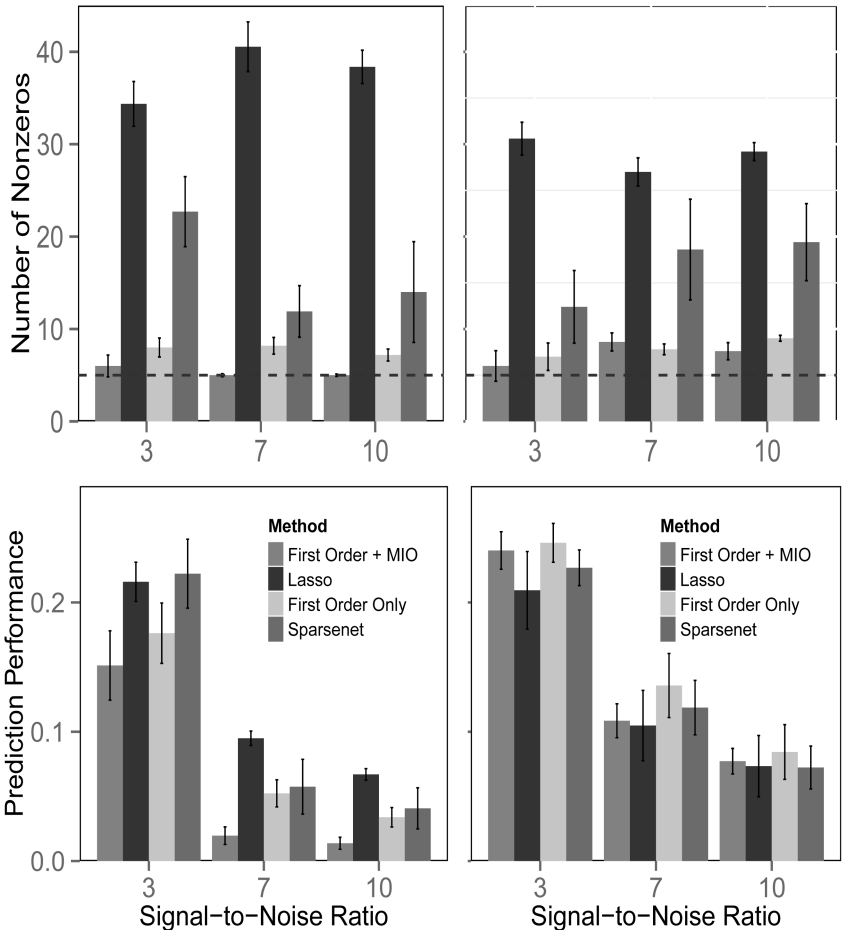


Figure 9: The sparsity and predictive performance for different procedures: [Left Panel] shows Example 1 with $n = 50, p = 1000, \rho = 0.8, k_0 = 5$ and [Right Panel] shows Example 2 with $n = 30, p = 1000$ —for each instance several SNR values have been shown.

coefficients in the model (with shrinkage). MIO (with warm-start) performs the best—both in terms of predictive accuracy and in selecting a sparse set of coefficients. MIO obtains the sparsest model among the four methods and seems to find better solutions in terms of statistical properties than the models obtained by the first order methods alone. Interestingly, the “optimal model” selected by the first order methods is more dense than that selected by the MIO. The number of non-zero coefficients selected by MIO remains fairly stable across different SNR values, unlike the other three methods.

In Figure 9 the right panel shows Example 2, with $n = 30, p = 1000, k_0 = 5$ and all non-zero coefficients equal one. In this example, all the methods perform similarly in terms of predictive accuracy. This is because all non-zero

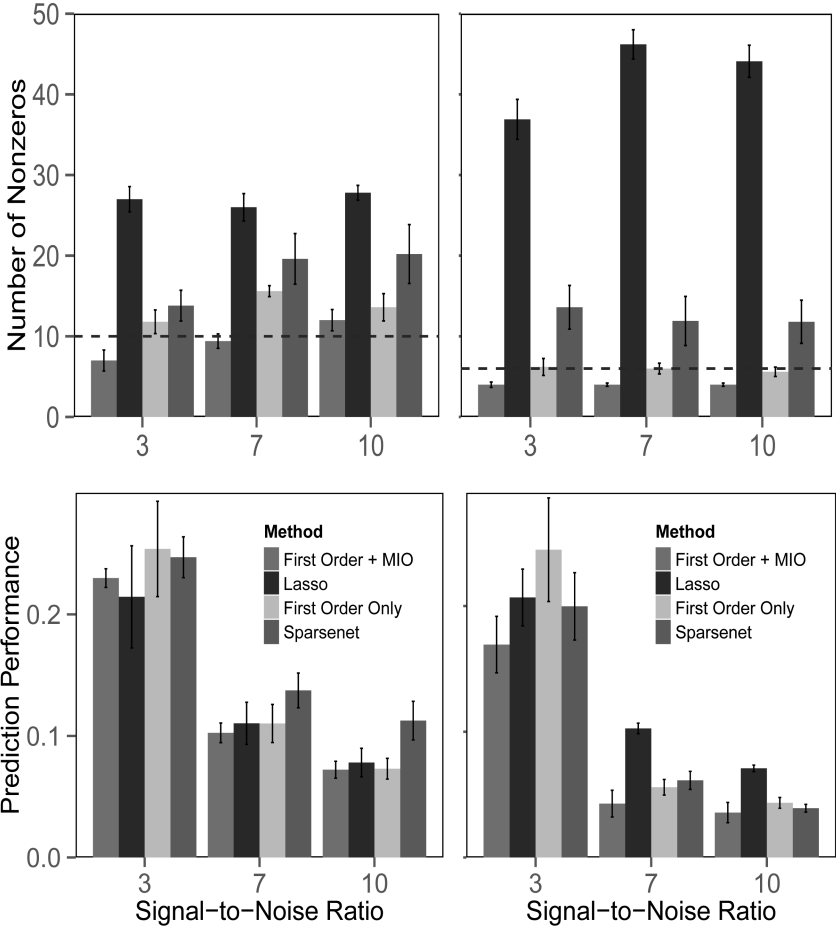


Figure 10: [Left Panel] Shows performance for data generated according to Example 3 with $n = 30, p = 1000$ and [Right Panel] shows Example 4 with $n = 50, p = 2000$.

coefficients in β^0 have the same value. In fact for the smallest value of SNR, the Lasso achieves the best predictive model. In all the cases however, the MIO achieves the sparsest model with favorable predictive accuracy.

In Figure 10, for both the examples, the model matrix is an iid Gaussian ensemble. The underlying regression coefficient β^0 however, is structurally different than Example 2 (as in Figure 9, right-panel). The structure in β^0 is responsible for different statistical behaviors of the four methods across Figures 9 (right-panel) and Figure 10 (both panels). The alternating signs and varying amplitudes of β^0 are responsible for the poor behavior of Lasso. The MIO (with warm-starts) seems to be the best among all the methods. For Example 3 (Figure 10, left panel) the predictive performances of Lasso and MIO are comparable—the MIO however delivers much sparser models than the Lasso.

The key conclusions are as follows:

1. The MIO best subset algorithm has a significant edge in detecting the correct sparsity structure for all examples compared to Lasso, Sparsenet and the stand-alone discrete first order method.
2. For data generated as per Example 1 with large values of ρ , the MIO best subset algorithm gives better predictive performance compared to its competitors.
3. For data generated as per Examples 2 and 3, MIO delivers similar predictive models like the Lasso, but produces much sparser models. In fact, Lasso seems to perform marginally better than MIO, as a predictive model for small values of SNR.
4. For Example 4, MIO performs the best both in terms of predictive accuracy and delivering sparse models.

5 Computational Results for Subset Selection with Least Absolute Deviation Loss

In this section, we demonstrate how our method can be used for the best subset selection problem with LAD objective (47).

Since the main focus of this paper is the least squares loss function, we consider only a few representative examples for the LAD case. The LAD loss is appropriate when the error follows a heavy tailed distribution. The datasets used for the experiments parallel those described in Section 4.1, the difference being in the distribution of ϵ . We took ϵ_i iid from a double exponential distribution with variance σ^2 . The value of σ^2 was adjusted to get different values of SNR.

Datasets analysed We consider a set-up similar to Example 1 (Section 4.1) with $k_0 = 5$ and $\rho = 0.9$. Different choices of (n, p) were taken to cover both the overdetermined ($n = 500, p = 100$) and high-dimensional cases ($n = 50, p = 1000$ and $n = 500, p = 1000$).

The other competing methods used for comparison were (a) discrete first order method (Section (3.4)) (b) MIO warm-started with the first order solutions and (c) the LAD loss with ℓ_1 regularization:

$$\min \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_1 + \lambda\|\boldsymbol{\beta}\|_1,$$

which we denote by LAD-Lasso. The training, validation and testing were done in the same fashion as in the least squares case. For each method, we report the number of non-zeros in the optimal model and associated prediction accuracy (53).

Figure 11 compares the MIO approach with others for LAD in the overdetermined case ($n > p$). Figure 12 does the same for the high-dimensional case ($p \gg n$). The conclusions parallel those for the least squares case. Since, in

the example considered, the features corresponding to the non-zero coefficients are weakly correlated and a feature having a non-zero coefficient is highly correlated with a feature having a zero coefficient—the LAD-Lasso selects an overly dense model and misses out in terms of prediction error. Both the MIO (with warm-starts) and the discrete first order methods behave similarly—much better than ℓ_1 regularization schemes. As expected, we observed that subset selection with least squares loss leads to inferior models for these examples, due to a heavy-tailed distribution of the errors.

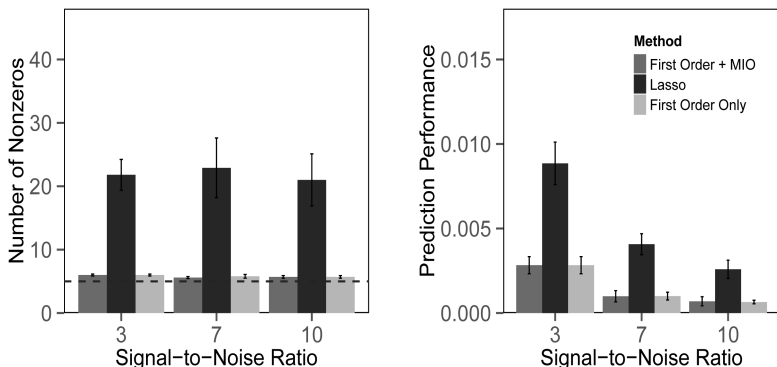


Figure 11: The sparsity and predictive performance for different procedures for $n = 500, p = 100$ for Problem (47). The data is generated as per Example 1 with $\rho = 0.9, k_0 = 5$ and double exponential errors—further details are available in the text. The acronym “Lasso” refers to LAD-Lasso (5). The MIO is seen to deliver sparser models with better predictive accuracy when compared to the LAD-Lasso.

The results in this section are similar to the least squares case. The MIO approach provides an edge both in terms of sparsity and predictive accuracy compared to Lasso both for the overdetermined and the high-dimensional case.

6 Conclusions

In this paper, we have revisited the classical best subset selection problem of choosing k out of p features in linear regression given n observations using a modern optimization lens, i.e., MIO and a discrete extension of first order methods from continuous optimization. Exploiting the astonishing progress of MIO solvers in the last twenty-five years, we have shown that this approach solves problems with n in the 1000s and p in the 100s in minutes to provable optimality, and finds near optimal solutions for n in the 100s and p in the 1000s in minutes. Importantly, the solutions provided by the MIO approach significantly outperform other state of the art methods like Lasso in achieving sparse models with good predictive power. Unlike all other methods, the MIO approach always provides a guarantee on its sub-optimality even if the algorithm is terminated early. Moreover, it can accommodate side constraints

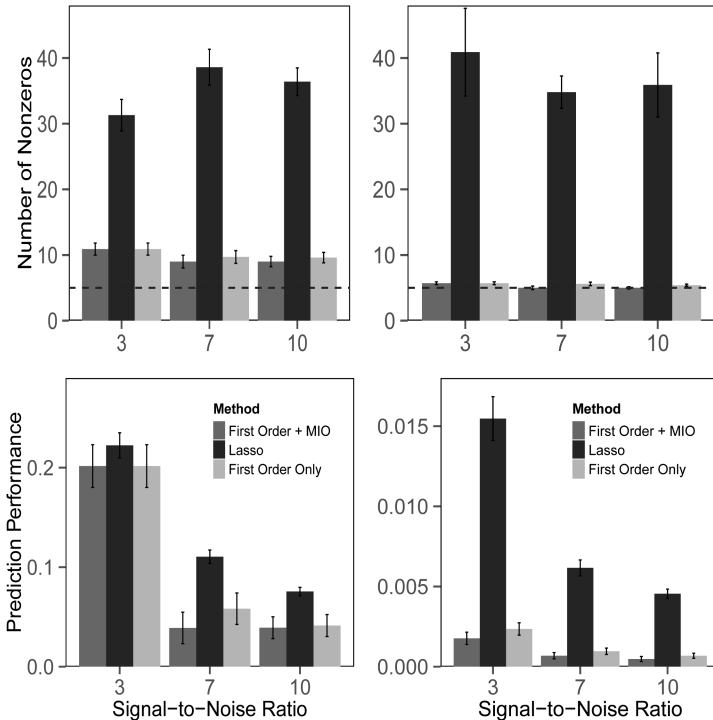


Figure 12: Figure showing the number of nonzero values and predictive performance for different values of n and p for Problem (47) (as in Figure 11). [Left panel] has $n = 50, p = 1000$ and [Right panel] has $n = 500, p = 1000$.

on the coefficients of the linear regression and also extends to finding best subset solutions for the least absolute deviation loss function.

While continuous optimization methods have played and continue to play an important role in statistics over the years, discrete optimization methods have not. The evidence in this paper as well as in [2] suggests that MIO methods are tractable and lead to desirable properties (improved accuracy and sparsity among others) at the expense of higher, but still reasonable, computational times.

References

- [1] Top500 Supercomputer Sites, Directory page for Top500 lists. Result for each list since June 1993. <http://www.top500.org/statistics/sublist/>. Accessed: 2013-12-04.
- [2] D. Bertsimas and R. Mazumder. Least quantile of squares regression via modern optimization. *Annals of Statistics*, 2014, to appear.
- [3] D. Bertsimas and R. Shioda. Algorithm for cardinality-constrained quadratic optimization. *Computational Optimization and Applications*, 43(1):1–22, 2009.
- [4] D. Bertsimas and R. Weismantel. *Optimization over integers*. Dynamic Ideas Belmont, 2005.
- [5] P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, pages 1705–1732, 2009.

- [6] D. Bienstock. Computational study of a family of mixed-integer quadratic programming problems. *Mathematical programming*, 74(2):121–140, 1996.
- [7] R. E. Bixby. A brief history of linear and mixed-integer programming computation. *Documenta Mathematica, Extra Volume: Optimization Stories*, pages 107–121, 2012.
- [8] T. Blumensath and M. Davies. Iterative thresholding for sparse approximations. *Journal of Fourier Analysis and Applications*, 14(5-6):629–654, 2008.
- [9] T. Blumensath and M. Davies. Iterative hard thresholding for compressed sensing. *Applied and Computational Harmonic Analysis*, 27(3):265–274, 2009.
- [10] P. Bühlmann and S. van-de-Geer. *Statistics for high-dimensional data*. Springer, 2011.
- [11] E. Candes, M. Wakin, and S. Boyd. Enhancing sparsity by reweighted ℓ_1 minimization. *Journal of Fourier Analysis and Applications*, 14(5):877–905, 2008.
- [12] E. Candes. The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathématique*, 346(9):589–592, 2008.
- [13] E. Candès and Y. Plan. Near-ideal model selection by ℓ_1 minimization. *The Annals of Statistics*, 37(5A):2145–2177, 2009.
- [14] E. Candes and T. Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *Information Theory, IEEE Transactions on*, 52(12):5406–5425, 2006.
- [15] S. Chen, D. Donoho, and M. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.
- [16] M. Dettling. Bagboosting for tumor classification with gene expression data. *Bioinformatics*, 20(18):3583–3593, 2004.
- [17] D. Donoho. For most large underdetermined systems of equations, the minimal ℓ^1 -norm solution is the sparsest solution. *Communications on Pure and Applied Mathematics*, 59:797–829, 2006.
- [18] D. Donoho and I. Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81:425–455, 1993.
- [19] D. Donoho and M. Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–2202, 2003.
- [20] D. Donoho and X. Huo. Uncertainty principles and ideal atomic decomposition. *Information Theory, IEEE Transactions on*, 47(7):2845–2862, 2001.
- [21] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression (with discussion). *Annals of Statistics*, 32(2):407–499, 2004. ISSN 0090-5364.
- [22] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360(13), 2001.
- [23] I. Frank and J. Friedman. A statistical view of some chemometrics regression tools (with discussion). *Technometrics*, 35(2):109–148, 1993.
- [24] J. Friedman. Fast sparse regression and classification. Technical report, Department of Statistics, Stanford University, 2008.
- [25] J. Friedman, T. Hastie, H. Hoefling, and R. Tibshirani. Pathwise coordinate optimization. *Annals of Applied Statistics*, 2(1):302–332, 2007.
- [26] G. Furnival and R. Wilson. Regression by leaps and bounds. *Technometrics*, 16: 499–511, 1974.
- [27] E. Greenshtein. Best subset selection, persistence in high-dimensional statistical learning and optimization under ℓ_1 constraint. *The Annals of Statistics*, 34(5):2367–2386, 2006.
- [28] E. Greenshtein and Y. Ritov. Persistence in high-dimensional linear predictor selection and the virtue of overparametrization. *Bernoulli*, 10(6):971–988, 2004.
- [29] I. Gurobi Optimization. Gurobi optimizer reference manual, 2013. URL <http://www.gurobi.com>.
- [30] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning, Second Edition: Data Mining, Inference, and Prediction (Springer Series in Statistics)*. Springer New York, 2 edition, 2009. ISBN 0387848576.
- [31] K. Knight and W. Fu. Asymptotics for lasso-type estimators. *Annals of Statistics*, 28(5):1356–1378, 2000.
- [32] P.-L. Loh and M. Wainwright. Regularized m-estimators with nonconvexity: Statistical

and algorithmic theory for local optima. In *Advances in Neural Information Processing Systems*, pages 476–484, 2013.

- [33] R. Mazumder, J. Friedman, and T. Hastie. Sparsenet: Coordinate descent with non-convex penalties. *Journal of the American Statistical Association*, 117(495):1125–1138, 2011.
- [34] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- [35] A. Miller. *Subset selection in regression*. CRC Press Washington, 2002.
- [36] B. Natarajan. Sparse approximate solutions to linear systems. *SIAM journal on computing*, 24(2):227–234, 1995.
- [37] G. Nemhauser. Integer programming: the global impact. Presented at EURO, INFORMS, Rome, Italy, 2013. http://euro2013.org/wp-content/uploads/Nemhauser_EuroXXVI.pdf. Accessed: 2013-12-04.
- [38] Y. Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming, Series A*, 103:127–152, 2005.
- [39] Y. Nesterov. Gradient methods for minimizing composite objective function. Technical report, Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, 2007. Technical Report number 76.
- [40] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer, Norwell, 2004.
- [41] G. Raskutti, M. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over-balls. *Information Theory, IEEE Transactions on*, 57(10):6976–6994, 2011.
- [42] R. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1996. ISBN 0691015864. URL <http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/0691015864>.
- [43] X. Shen, W. Pan, Y. Zhu, and H. Zhou. On constrained and regularized high-dimensional regression. *Annals of the Institute of Statistical Mathematics*, 65(5): 807–832, 2013.
- [44] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- [45] R. Tibshirani. Regression shrinkage and selection via the lasso: a retrospective. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(3):273–282, 2011.
- [46] J. Tropp. Just relax: Convex programming methods for identifying sparse signals in noise. *Information Theory, IEEE Transactions on*, 52(3):1030–1051, 2006.
- [47] M. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using constrained quadratic programming (lasso). *Information Theory, IEEE Transactions on*, 55(5):2183–2202, 2009.
- [48] C.-H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38(2):894–942, 2010.
- [49] C.-H. Zhang and J. Huang. The sparsity and bias of the lasso selection in high-dimensional linear regression. *Annals of Statistics*, 36(4):1567–1594, 2008.
- [50] C.-H. Zhang and T. Zhang. A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical Science*, 27(4):576–593, 2012.
- [51] T. Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *The Journal of Machine Learning Research*, 11:1081–1107, 2010.
- [52] Y. Zhang, M. Wainwright, and M. I. Jordan. Lower bounds on the performance of polynomial-time algorithms for sparse linear regression. *arXiv preprint arXiv:1402.1918*, 2014.
- [53] P. Zhao and B. Yu. On model selection consistency of lasso. *Journal of Machine Learning Research*, 7:2541–2563, 2006.
- [54] H. Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101:1418–1429, 2006.
- [55] H. Zou and R. Li. One - step sparse estimates in nonconcave penalized likelihood problems. *The Annals of Statistics*, 36(4):1509–1533, 2008.



DIMITRIS BERTSIMAS

DIMITRIS BERTSIMAS

Dimitris Bertsimas is the Boeing Professor of operations research and the co-director of the Operations Research Center at MIT. He received his SM and PhD in applied mathematics and operations research from MIT in 1987 and 1988 respectively. He has been with the MIT faculty since 1988. His research interests include optimization, statistics, and applied probability and their applications in health care, finance, operations management, and transportation. He has co-authored more than 150 scientific papers and three graduate level textbooks. He is currently department editor in optimization for *Management Science* and former area editor in operations research in financial engineering. He has supervised over 70 doctoral students. He is a member of the National Academy of Engineering, and has received numerous research awards including the Morse Award (2013), Pierskalla Award (2013), TSL Best Paper Award (2013), Farkas Prize (2008), Erlang Prize (1996), SIAM Activity Group on Optimization (SIAG/OPT) Prize (1996), Bodossaki Prize (1998), and the Presidential Young Investigator Award (1991–1996). Bertsimas has co-founded several companies in the areas of financial services, health care, aviation, and publishing.



PHILIP McCORD MORSE

Philip McCord Morse⁴

1903–1985

by

Robert Herman

Philip McCord Morse, professor emeritus of physics at the Massachusetts Institute of Technology (MIT), founder and pioneer of modern operations research, physicist and Renaissance scientist, community leader, and leader in professional societies, died on September 5, 1985, in Concord, Massachusetts. As Phil Morse wrote in his autobiography, *In at the Beginnings: A Physicist's Life*, 1977, "They told me I was born on August 6, 1903, at three in the morning; I don't remember. My seventy-year memory tape is a series of vividly recollected scenes, separated by blanks later filled in with conjecture and hearsay. The early scenes are disconnected flashes, glimpses of a now unfamiliar world, seen through a stranger's eyes. It takes effort to remember how different that world was, how many differences there are between the Midwest of 1910 and the East Coast of the 1970s."

Morse's distinguished career in science and technology is characterized by a remarkable breadth and diversity of interests. In physics, it ranged from acoustics and quantum mechanics to nuclear physics and methods of theoretical physics. In operations research, which he pioneered, his career encompassed military operations research, vehicular traffic, queues, and public systems. His fundamental contributions in these diverse areas, together with his service to the professional community and society in general, created a most outstanding career.

His early developmental years were spent in Cleveland, Ohio. He was the son of a telephone engineer, the grandson of a civil engineer, and the great grandson of an architect and builder. His great grandfather worked for the federal government designing and building post offices and custom houses all over the country and was also elected to the Ohio legislature. While still in grade school, Morse read voraciously, was attracted to chemistry, and learned to play the violin. He indicated that while facts didn't interest him very much he was excited by patterns, such as the recurrent patterns in the Mendeleev table of the elements. During high school he decided to become a chemist. Interestingly, he never aspired to be a mathematician because, he said, mathematics had been treated as a tool rather than as a subject for intellectual exploration. Eric Bell's *Men of Mathematics* had not yet been written when Morse made that statement; he later speculated that if the book

⁴This essay, reprinted with permission of the National Academy of Engineering, appeared in Memorial Tributes, National Academy of Engineering, Volume 4, National Academy Press, Washington, D.C., 1991.

had already appeared he might have become enmeshed in the mysteries of prime numbers of Diophantine analysis and his entire life might have been different. As for his nonscholastic interests, when the radio craze hit Cleveland in the early twenties, Morse operated his own radio supply and repair shop.

After one year of undergraduate study, Morse took the year 1922-23 off to operate his radio business when family fortunes were at a low ebb. By the fall of 1923 when he returned to college as a sophomore, he was considerably more certain about what he wished to learn. Upon deciding to pursue the physics program, his father's only comment was, "That's fine, but what will you do for money?" It is interesting to read in Morse's recollections that he didn't share this concern for money and that he envisioned a career teaching college physics. He commented at the time that, "Professors never got rich-but then they never seemed to starve."

Morse received his B.S. in 1926 from what was then the Case School of Applied Science. He pursued his graduate studies at Princeton University and received his Ph.D. in physics in 1929. It was during his undergraduate days that he became involved with the eminent American physicist Dayton C. Miller, who was one of the earliest experts in sound and musical acoustics, and whose large collection of flutes is now in the Library of Congress. It was during this period that Morse developed his lifelong interest in acoustics.

Physics and mathematics claimed much of his time as a graduate student at Princeton. Three courses didn't sound like much to him, but analytic dynamics, electron theory and mathematical physics generated a great work load. Unlike the students of pure mathematics, Morse was interested in analysis and higher algebra as the language of physics. The late 1920s were exciting times thanks to the development of the new quantum mechanics; in 1930 Dirac prophesied accurately that quantum mechanics would explain all of chemistry and most of physics.

Aside from his course work and research on molecular physics with Ernst Stueckelberg, with whom he published several papers, Morse developed a solution for a force that was repulsive when two particles are close together, attractive when they are further apart, and under which they vanish at a greater distances. He realized that he had stumbled upon a quantum mechanical representation of a vibrating diatomic molecule. To this day, the particular force field, expressed as a related potential field, is known as the Morse Potential.

Edward Condon, upon his return from Europe, where the new quantum theory had been developed, decided to write an English text on the subject. When the writing progressed too slowly, he invited Morse to collaborate. The idea appealed to Morse as an opportunity to learn the rapidly developing quantum mechanics not only by teaching it but by structuring a monograph on it. Thus, Morse coauthored one of the earliest texts on the new quantum theory.

Among his other notable associations, he assisted in the development of the theoretical understanding of the Davison-Germer experiment during a summer at the AT&T Bell Laboratories. His postdoctoral studies were conducted with Arnold Sommerfeld in Munich and included theoretical research in electron scattering under an international fellowship. Thanks to Morse's early renown, Karl T. Compton, then president of MIT, asked Morse to join the MIT physics

faculty when he returned from his fellowship in Europe. As Morse recounts, "It was easy to say yes."

Morse joined the MIT physics faculty in 1931 as assistant professor, rapidly rose to associate professor in 1934, and became a full professor in 1938. With his very broadly gauged interests, he participated in the development of the physics curriculum and accepted the position of graduate registration officer. His research continued in a diverse fashion; during this period he worked on electron scattering, nuclear binding forces, and even on the subject of stellar interiors in astrophysics. One of his important contributions to physics was the acoustics textbooks *Vibrations and Sound* published in 1936. This work presented the application of scattering theory to sound waves. In fact, it was also during this early period in Morse's life that he developed course notes that were later combined with those of Herman Feshbach to produce the famous two-volume work *Methods of Theoretical Physics*, published in 1953. The book is a basic source of methods of mathematical physics to this day.

With the advent of World War II, Philip Morse's Renaissance talent entered a new phase in his technical life. By the time the United States entered the war, the catastrophic loss of allied ships to the German U-boats in the Atlantic Ocean was a major concern. It was imperative that the U.S. develop superior equipment that would locate and neutralize this threat. The British, who had been engaged in the struggle for two years, already had several operations research groups not only designing equipment but also studying and maximizing its effectiveness in actual war operations. Early in 1942 the U.S. armed forces established an operations research group in the navy. Morse, who was considered a distinguished scientist and who had been the director of a project at the Underwater Sound Laboratory at Harvard University for the previous two years, was chosen by the National Defense Research Council to head the operations research effort.

Several months after the formation of the operations research group, the navy consolidated the antisubmarine operations under the Tenth Fleet, and the Antisubmarine Warfare Operations Research Group was transferred to Washington, DC. Morse had a substantial fraction of the group out in the field working with the operational commands. He did an outstanding job both in coordinating the technical work and in his liaison with the operational leaders running the actual war operations. Those who worked with Morse during this period report that it was a continuous learning experience. As the war effort and operations research became more successful, the Antisubmarine Warfare Operations Research Group became the navy's Operations Research Group. This group took on submarine activity studies in the Pacific Theater of Operations. It then addressed naval air activities and ultimately became involved in all aspects of navy task force operations. The group became very well accepted and at the conclusion of the war Morse received the Presidential Medal for Merit, the nation's highest civilian award.

After the war Morse generated an orderly windup of the group's activities, part of which became the nucleus of the Operations Evaluation Group. He returned to research and teaching at MIT but continued to monitor this postwar transition. In 1946, he had been at MIT no longer than one year when he became the director of the Atomic Energy Commission's Brookhaven

National Laboratory. The position occupied all of his time in organization and administration and left no time for personal scientific research. In 1948, with Brookhaven well established, Morse went to Washington to organize an operations research group for the Secretary of Defense and the Joint Chiefs of Staff. This resulted in the Weapons Systems Evaluation Group for which he served as deputy director and director of research until 1950. The group's civilian unit developed into the Institute for Defense Analyses in 1956; Morse served as a trustee.

In another area of interest, Morse was convinced of the great importance of computation and the rapidly growing power of the digital computer. This no doubt arose from his experience with calculations in acoustics and astrophysics in the late 1930s. The establishment of the MIT Computation Center was a result of his efforts to introduce computers to research and research to computers in the late 1940s and early 1950s. He became its first director and served in that position until 1967.

In 1952 Morse created an operations research activity at MIT with an interdepartmental committee and a small contract for fundamental research from the U.S. Army. In two years, the first doctoral student, John D.C. Little, was graduated and in 1956 the Operations Research Center was formally established with Morse as director; he remained in this role until his official retirement in 1968. His high research activity in the field of operations research was continuous and included the following books: *Queues, Inventories and Maintenance*, 1958; *Library Effectiveness: A Systems Approach*, 1968; *Operations Research for Public Systems*, 1967, coeditor; and *Analysis of Public Systems*, 1972, coeditor.

The Operations Research Society of America was founded in 1952, and as might have been expected, Morse became its first president. Of the next eight presidents, half had worked for him in one capacity or another, mostly during World War II. About twenty years later, there came an echo of Morse's influence as two of his former students became presidents of the society. Morse received the Frederick W. Lanchester Prize of the Operations Research Society in 1968 for his library work and was the first recipient of that society's George E. Kimball Medal in 1974 for his contributions to the profession of operations research in general and to the society in particular.

Professor Morse's worldwide promotion of operations research never ceased. He was involved in organizing the first International Operations Research Conference in 1957; the International Federation of Operations Research Societies originated at this conference. Interest in the operations research discipline overseas led to the 1959 NATO conference with Morse as chairman of the advisory panel. He was associated with many international operations research projects in which he always stressed that the discipline was applicable to a host of fundamental problems that were neither military nor industrial in nature. It is interesting to recall that most recently, in April 1985, at the age of 81, Morse chaired a session at ORSA's Boston meeting and spoke on the early use of computers in operations research, a topic that combined two of his major interests.

Morse's honors are legion. Among these, he was a member of the National Academy of Sciences; and a fellow of the American Academy of Arts and

Sciences, the Acoustical Society of America, and the American Physical Society. He was elected to the National Academy of Engineering in 1985. He was also a member of Sigma Xi, Tau Beta Pi, and the Cosmos Club of Washington. He received the Silver Medal of the Operational Research Society of the United Kingdom, and the Gold Medal of the Acoustical Society of America. He was the president of the Acoustical Society of America (1974–1977) and chairman (1975) of the Governing Board of the American Institute of Physics. From 1958 to 1960 he was chairman of the MIT faculty.

Philip Morse, one of the first wave of home-grown American scientists, made outstanding contributions to science and technology through his work in physics, computer science and operations research. He influenced and guided many students and colleagues in the struggle to seek scientific truth. In his autobiography Morse gives great food for thought to many of us. He reflects that his successes would have been fewer had he not chosen, back in 1923, to become a physicist through training that forced him to look facts in the face, that made him want to measure them and work out their implications, whether these facts applied to atoms or automobiles.

The last comment of Morse's autobiography conveys much of his philosophy: "For those who like exploration, immersion in scientific research is not dehumanizing; in fact, it is a lot of fun. And, in the end, if one is willing to grasp the opportunities, it can enable one to contribute something to human welfare."



*Institute for Operations Research and the Management Sciences
5521 Research Park Drive, Suite 200
Catonsville, MD 21228*