

ARTICLE

Sparse Convex Optimization Toolkit: A Mixed-Integer FrameworkAlireza Olama^a, Eduardo Camponogara^a, and Jan Kronqvist^b^aAutomation and System Engineering Department, Federal University of Santa Catarina, Florianópolis, Brazil^bDepartment of Mathematics, KTH Royal Institute of Technology, Stockholm, Sweden**ARTICLE HISTORY**

Compiled November 1, 2022

ABSTRACT

This paper proposes an open-source distributed solver for solving Sparse Convex Optimization (SCO) problems over computational networks. Motivated by past algorithmic advances in mixed-integer optimization, the Sparse Convex Optimization Toolkit (SCOT) adopts a mixed-integer approach to find exact solutions to SCO problems. In particular, SCOT brings together various techniques to transform the original SCO problem into an equivalent convex Mixed-Integer Nonlinear Programming (MINLP) problem that can benefit from high-performance and parallel computing platforms. To solve the equivalent mixed-integer problem, we present the Distributed Hybrid Outer Approximation (DiHOA) algorithm that builds upon the LP/NLP based branch-and-bound and is tailored for this specific problem structure. The DiHOA algorithm combines the so-called single- and multi-tree outer approximation, naturally integrates a decentralized algorithm for distributed convex nonlinear subproblems, and utilizes enhancement techniques such as quadratic cuts. Finally, we present detailed computational experiments that show the benefit of our solver through numerical benchmarks on 140 SCO problems with distributed datasets. To show the overall efficiency of SCOT we also provide performance profiles comparing SCOT to other state-of-the-art MINLP solvers.

KEYWORDS

sparse optimization; mixed integer nonlinear programming; distributed computing, outer approximation

1. Introduction

In recent years, *Sparse Convex Optimization (SCO)* has gained considerable attention in several disciplines, from machine learning and engineering to economics and finance [4, 6, 39]. Several mathematical optimization problems in this context can be formulated as a general convex optimization problem subject to a constraint that allows only up to a certain number of decision variables to be nonzero. We refer to this constraint as a *sparsity constraint*. Hence, any convex optimization problem with the sparsity constraint can be regarded as a SCO problem [2, 6, 9, 36].

Due to the non-convexity and discontinuity of the sparsity constraint, the SCO problems are known to be NP-Hard [6, 27]. To overcome the computational difficulties imposed by the sparsity constraint, computationally tractable convex optimization-based methods have been proposed. One of the popular methods is a ℓ_1 norm relaxation

Table 1. Acronyms

Acronym	Meaning
API	Application Programming Interface
BnB	Branch and Bound
CLI	Command Line Interface
DiHOA	Distributed Hybrid Outer Approximation
DiPOA	Distributed Primal Outer Approximation
D-MINLP	Distributed Mixed Integer Nonlinear Programming
DSLInR	Distributed Sparse Linear Regression
DSLogR	Distributed Sparse Logistic Regression
LFC	Local Fusion Center
MILP	Mixed Integer Linear Programming
MINLP	Mixed Integer Nonlinear Programming
MIP	Mixed Integer Programming
MPI	Message Passing Interface
NLP	Nonlinear Programming
OA	Outer Approximation
QCLP	Quadratically Constrained Linear Programming
RH-ADMM	Relaxed Hybrid Alternating Direction Method of Multipliers
SCO	Sparse Convex Optimization
SCOT	Sparse Convex Optimization Toolkit
SOS-1	Specially Ordered Set of Type I

method where a ℓ_1 regularizer is imposed on the decision vector. The ℓ_1 method naturally produces a sparse solution by setting many variables to zero. One of the popular ℓ_1 based methods is Lasso which is widely used in statistics and machine learning communities [10, 11, 38].

One of the important reasons behind the popularity of ℓ_1 based methods is their computational efficiency and scalability to practical-sized problems. However, in spite of their favorable computational properties, these methods can have some shortcomings. For example, they cannot guarantee that ℓ_1 based methods find the correct sparsity for general problems. Moreover, in some applications, the desired sparsity structure is different from general sparsity and cannot be easily obtained by a ℓ_1 regularization. An example of such a sparse structure is *group sparsity* in which a block or a group of independent variables are either all zero or all nonzero. Some notable applications with group sparsity are block-wise linear regression [20], logistic regression [4], compressed sensing [13], and microarray analysis [26].

Another approach to solving the SCO problems is to view the SCO problems as equivalent Mixed Integer Programming (MIP) problems. Considering recent advances in mixed-integer optimization algorithms and technologies, the MIP problems can be solved efficiently by current mixed-integer optimization solvers such as Gurobi [18]. The resulting MIP formulation is flexible and can be adjusted based on the application’s needs. Moreover, by defining suitable binary variables, the MIP framework can provide exact sparse solutions to SCO problems. The MIP framework is gaining popularity in various areas such as statistical data analysis and interpretable machine learning [3, 5–7], sparse control [1], unit commitment, face recognition, sensor network design [24], portfolio optimization [2], and compressed sensing [15].

The mentioned works focused on *centralized* solutions to SCO problems that might not be suitable for modern real-world applications when the data is inherently distributed or available in large volumes. Considering the limitations of centralized architectures, in the past few years, mainly because of the rise of Big data, distributed optimization over networks has gained growing attention [30, 31]. The primary purpose of distributed optimization is to solve an optimization problem over a network of computing nodes. Each node performs local computations with access only to a

portion of the problem data. Nodes are also capable of exchanging information with other nodes in the network. See [31] for a comprehensive overview of the most common distributed optimization algorithms.

This paper introduces **SCOT**, a distributed optimization solver designed to solve SCO problems over peer-to-peer networks of computing nodes. In essence, **SCOT** consists of two distributed algorithms developed by the authors to solve SCO problems where an iterative procedure is applied by each network node, alternating communication, and computation phases until a solution is found. To the best of our knowledge, **SCOT** is the first software framework that can solve SCO problems using distributed algorithms, enabling practical applications with a large number of sample data points, while keeping the data private to each node and allowing implementation in a computer cluster.

Formally, we consider the SCO problem as a mathematical programming problem that consists of finding the κ -sparse optimal solution of a distributed convex optimization problem of the following form,

$$\min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^N f_i(\mathbf{x}) \quad (1a)$$

$$\text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b} \quad (1b)$$

$$\|\mathbf{x}\|_0 \leq \kappa \quad (1c)$$

where N is the number of nodes of the computation network, $\mathbf{x} \in \mathbb{R}^n$ is the vector of decision variables, and $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function assumed to be continuously differentiable and only known by node i , for all $i \in \{1, \dots, N\}$. We use the ℓ_0 norm (*i.e.*, $\|\mathbf{x}\|_0 = |\text{supp}(\mathbf{x})| = |\{j : x_j \neq 0\}|$) to define the sparsity constraint, which imposes the number of non-zero elements of \mathbf{x} to be less than a given integer κ . Finally, the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and the vector $\mathbf{b} \in \mathbb{R}^m$ define the set of given linear constraints assumed to be known by all nodes.

1.1. Main Contributions

The main idea behind **SCOT** is to transform problem (1) into an equivalent *Distributed Mixed Integer Nonlinear Programming (D-MINLP)* problem. By using the concept of the so-called Local Fusion Centers (LFCs) presented in [32], we recast problem (1) as a consensus optimization problem and introduce several constraints to model the sparsity constraint (1c). Based on the Outer Approximation (OA) algorithm [12, 17, 21], **SCOT** consists of two main algorithms, namely, *Distributed Primal Outer Approximation (DiPOA)* and *Distributed Hybrid Outer Approximation (DiHOA)*.

The DiPOA algorithm proposed in [33] extends the OA algorithm by embedding a fully decentralized algorithm, namely the Relaxed Hybrid Alternating Direction Method of Multipliers (RH-ADMM) [32]. The RH-ADMM assumes a particular hybrid architecture on the computational network, developed to solve distributed convex optimization problems. In particular, DiPOA solves the primal problem of the OA using the RH-ADMM algorithm and handles distributedly the demanding computational part of the OA algorithm that deals with the solution of convex NLPs. In practice, there exist many cases where a significant part of the solution time is spent on solving the NLP problems. For example, in Table 1 [21], it can be seen that OA spends more than 150 seconds on solving NLP problems when solving a moderate-size convex MINLP. Moreover, for inherently distributed problems in which the data is spread

over a possibly large computational network, a single NLP will be challenging, if even possible, to solve in a classical centralized fashion. For such inherently decentralized problems, a distributed algorithm can offer great computational advantages.

Despite solving problem (1) distributedly, DiPOA is developed based on a multiple-tree OA algorithm whereby a BnB tree is built from scratch at each iteration of the algorithm. Therefore, most of the overall solution time of DiPOA is usually spent on solving MIP sub-problems. In this paper, we tackle the limitations of DiPOA by proposing DiHOA, a distributed algorithm that improves DiPOA performance by gradually building a single BnB tree to avoid constructing and solving many similar MILP problems from scratch. The main idea of DiHOA is to solve problem (1) by dynamically updating the MILP subproblem, in a fashion similar to the LP/NLP-BnB presented by [34]. In principle, DiHOA starts with the multiple-tree search strategy up to a certain iteration and introduces high-quality cuts until a suitable event is triggered. Once the event is triggered, DiHOA switches from the multiple-tree search strategy to the single-tree search strategy that builds a single BnB tree, starting with multiple cuts initially introduced to its root node to augment the initial formulation. The single BnB tree then tightens up the integer relaxations by dynamically introducing more linear approximations (cuts) to the MILP problem. The multiple-tree search strategy is only applied in some initial iterations of the DiHOA since the MIP problems are typically much easier to solve, as the nonlinear constraints are only roughly represented through a few constraints. Moreover starting the BnB search with a tighter approximation of the nonlinear constraints can result in a much smaller BnB tree, as a smaller infeasible region of the continuously relaxed search space is explored. Without these initial cuts, the nonlinear constraints would be completely ignored until an integer solution is found and the nonlinear constraints would be poorly approximated until a few integer solutions is explored in the BnB tree.

The main contributions of this work are summarized as follows:

- SCOT, a distributed software framework to model and solve SCO problems.
- The distributed hybrid outer approximation algorithm to solve the SCO problem (1).
- Modeling and heuristic techniques to improve the efficiency of the proposed algorithm.
- A computational analysis of the proposed algorithms for various SCO real-world applications.

1.2. Paper Organization

The paper is organized as follows. In Section 2 we introduce the distributed SCO problem by reformulating problem (1) into an equivalent MINLP problem. Section 3 presents the primal and dual problems used by SCOT algorithms. The DiHOA algorithm is proposed in Section 4. Section 5 introduces SCOT and its main components. Finally, the numerical comparison and algorithm analysis are presented in Section 6. Section 7 concludes the paper.

2. Distributed Sparse Convex Optimization

Various modeling techniques are used by SCOT to find a solution to problem (1) efficiently. First, SCOT transforms problem (1) into a consensus optimization problem and

then multiple modeling techniques are used to handle the sparsity constraint (1c).

2.1. Consensus Optimization Modeling

As defined in problem (1), the decision variables are shared between the nodes. Each node only has information to construct its objective function, while keeping the local problem data private from other nodes. This distributed setting is a typical pattern in some significant learning and control applications. For instance, in distributed machine learning, an efficient technique to deal with large volumes of data consists of distributing the data over a network [28, 31]. In this case, a significant reduction in memory size for computation can be achieved while keeping the same unknown parameters of the model.

To decompose (1), we adhere to the concept of *hypergraphs*, which is a generalization of a regular graph in which an edge can join an arbitrary number of vertices. Multiple computational sources, called LFCs, can be employed by adopting this structure in the network. We consider a hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ defined as follows: $\mathcal{V} = \{1, 2, \dots, N\}$ is the set of nodes such that node i decides upon the values of vector variable \mathbf{x}_i ; $\mathcal{E} = \{\mathcal{E}_k \subset \mathcal{V} : k = 1, \dots, K\}$ is the set of hyperedges, where a hyperedge \mathcal{E}_k connects all nodes $i \in \mathcal{E}_k$ and K is the number of hyperedges. Now we introduce the concept of path in a hypergraph: a path $p(i, j) = \langle \mathcal{E}_1, \dots, \mathcal{E}'_k \rangle$ connects nodes i and j if $i \in \mathcal{E}'_1$, $j \in \mathcal{E}'_k$, and $\mathcal{E}'_l \cap \mathcal{E}'_{l+1} \neq \emptyset$ for $l = 1, \dots, k - 1$, and $\mathcal{E}'_l \in \mathcal{E}$ for all l . Put another way, through the hyperedges, a path $p(i, j)$ establishes a communication channel between nodes i and j .

We assume that the hypergraph \mathcal{H} is connected, meaning that for all $i, j \in \mathcal{V}$ there exists a path $p(i, j)$ connecting i and j . By using the hypergraph structure, the equivalent formulation for (1) is obtained as follows,

$$\min_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N \\ \mathbf{y}_1, \dots, \mathbf{y}_K}} \sum_{i=1}^N f_i(\mathbf{x}_i) \quad (2a)$$

$$\text{subject to } \mathbf{A}\mathbf{x}_i \leq \mathbf{b}, \forall i = 1, \dots, N, \quad (2b)$$

$$\mathbf{x}_i = \mathbf{y}_j, \forall i \in \mathcal{E}_j, \mathcal{E}_j \in \mathcal{E} \quad (2c)$$

$$\|\mathbf{y}_j\|_0 \leq \kappa, \forall j = 1, \dots, K \quad (2d)$$

where $\mathbf{x}_i \in \mathbb{R}^n$ are vectors of decision variables associated with the nodes and $\mathbf{y}_j \in \mathbb{R}^n$ are auxiliary variables associated with the LFCs, which are represented by the hyperedges.

2.2. Sparsity Constraint Modeling

Here, we present three modeling techniques implemented by SCOT to express the sparsity constraints (2d), namely the *big-M*, *Specially Ordered Set of Type I* (SOS-1), and a *hybrid* approach.

2.2.1. Big-M Method

The Big-M method is arguably the simplest technique for modeling the sparsity constraint. This method incorporates a binary variable and an estimated upper bound into the model for each continuous variable appearing in the sparsity constraint. By using

the Big-M method, we recast the sparsity constraint (2d) as the following inequalities,

$$-M_j\delta_{jk} \leq y_{jk} \leq M_j\delta_{jk}, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (3a)$$

$$\sum_{k=1}^n \delta_{jk} \leq \kappa, \forall j = 1, \dots, K \quad (3b)$$

$$\delta_{jk} \in \{0, 1\}, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (3c)$$

where y_{jk} is the k -th element of \mathbf{y}_j , $\boldsymbol{\delta}_j, \forall j = 1, \dots, K$, is a vector of binary variables whose k -th element is denoted by δ_{jk} , and M_j is a constant assumed to be a valid upper bound for $\|\mathbf{y}_j\|_\infty$. In this case, if $\delta_{jk} = 0$ then $y_{jk} = 0$ and $y_{jk} = 1$ otherwise. Thus, inequalities (3a) and (3b) impose the maximum number of nonzero variables in \mathbf{y}_j . The big-M parameter M_j is not known a priori, and a too small value of M_j may lead to a sub-optimal solution. A large M_j on the other hand, will result in a weak continuous relaxation and strongly affect the number of nodes that need to be explored in BnB. Hence, a good choice of M_j affects the strength of the formulation, being critical for MIP algorithms to obtain high-quality bounds. In the context of learning and control applications, the big-M value M_j can typically be computed from data in statistical learning tasks [7] and from the physical bounds in control applications [1].

2.2.2. Specially Ordered Set of Type I (SOS-1) Method

This section discusses the sparsity constraint reformulation using the SOS-1 constraint. Any feasible solution to problem (2) satisfies the following complementary constraints,

$$(1 - \delta_{jk})y_{jk} = 0, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (4)$$

which is equivalent to constraint (2d). In order for constraint (4) to be satisfied, either $(1 - \delta_{jk})$ or y_{jk} must be zero. Such constraints can be modeled via integer optimization software using Specially Ordered Sets of Type I (SOS-1) [8] as follows,

$$(y_{jk}, 1 - \delta_{jk}) : \text{SOS-1}, \forall j = 1, \dots, K, \forall k = 1, \dots, n. \quad (5)$$

It is worth noting that such a constraint is not implemented explicitly employing algebraic equations, as in the Big-M method. Instead, the optimization solver can directly enforces the SOS-1 constraint by branching on sets of variables. The MIP solvers usually adopt different approaches to handle the SOS-1 constraints. One approach is to use an equivalent big-M formulation in case strong bounds on variables can be obtained. Otherwise, the MIP solver can also handle the constraint directly through branching. By using the SOS-1 technique, we allow the MIP solver to choose the approach it determines to be the best choice to deal with the SOS-1 constraints.

2.2.3. MINLP Reformulation

This section discusses the primary problem formulation that SCOT considers. Although the big-M formulation is relatively straightforward, choosing a suitable big-M value M_j is important. In case M_j is too small, the optimization algorithm may cut off valid solutions. However, if M_j is excessively large, the model may become numerically difficult to solve. SOS-1 constraints, on the other hand, have the advantage of avoiding these

types of problems, as they do not rely on a problem-dependent constant value. For applications that require a relatively small value of M_j , the big-M modeling technique may be a proper choice, although, for other cases, the SOS-1 method may be suitable. In this regard, **SCOT** can provide the functionality to choose the most proper method depending on the application and problem data. Considering big-M and SOS-1 constraints for modeling the sparsity constraint, the main optimization problem that **SCOT** attempts to solve is expressed as follows,

$$\begin{aligned}
& \min_{\substack{\gamma \\ \mathbf{x}_1, \dots, \mathbf{x}_N \\ \mathbf{y}_1, \dots, \mathbf{y}_K}} \gamma & (6a) \\
\text{subject to } & \sum_{i=1}^N f_i(\mathbf{x}_i) - \gamma \leq 0 & (6b) \\
& \mathbf{A}\mathbf{x}_i \leq \mathbf{b}, \forall i = 1, \dots, N, & (6c) \\
& \mathbf{x}_i = \mathbf{y}_j, \forall i \in \mathcal{E}_j, \mathcal{E}_j \in \mathcal{E} & (6d) \\
& -M_j\delta_{jk} \leq y_{jk} \leq M_j\delta_{jk}, \forall j = 1, \dots, K, \forall k = 1, \dots, n & (6e) \\
& (y_{jk}, 1 - \delta_{jk}) : \text{SOS-1}, \forall j = 1, \dots, K, \forall k = 1, \dots, n & (6f) \\
& \sum_{k=1}^n \delta_{jk} \leq \kappa, \forall j = 1, \dots, K & (6g) \\
& \delta_{jk} \in \{0, 1\}, \forall j = 1, \dots, K, \forall k = 1, \dots, n & (6h)
\end{aligned}$$

in which the equivalent epigraph reformulation is used. The advantage of using both constraints (6e) and (6f) is that it might be possible to determine better M_j coefficients than the MIP solver is able to derive. In that way, **SCOT** provides more information to the MIP solver by including constraint (6e) and (6f). Problem (6) is an MINLP problem with a separable structure, where the objective function is linear and $\gamma \in \mathcal{R}$ is an auxiliary variable. In the following sections, we introduce a distributed formulation and algorithm that **SCOT** implements to solve problem (6).

3. SCOT Primal and Dual Problem

Here, we present the dual and primal problems and we discuss two algorithms that **SCOT** implements in the next section. Akin to other decomposition-based MINLP algorithms, **SCOT** decomposes problem (6) into two main sub-problems, namely, the *primal* and *dual* problems [21, 25]. We use the term *primal solution* and *primal bound* as the optimal solution and objective value of the primal problem respectively. Similarly, *dual solution* and *dual bound* are used for the dual problem. It should be noted that we use the terminology dual problem for a problem whose optimal solution provides a valid lower bound on the optimal objective value of problem (6) and whose feasible set contains all feasible solutions of problem (6).

The primal solution is assumed to satisfy all linear, nonlinear, and consensus constraints of problem (6) to a given tolerance. Additionally, the current best-known primal solution found by the algorithm is referred to as the *incumbent solution*, and its objective value is the current primal bound. The incumbent solution is updated when **SCOT** algorithms find a primal solution with a lower objective value. In the standard OA algorithm, the primal problem is a convex NLP problem; however, owing

to the distributed nature of problem (6), the primal problem of SCOT is a distributed convex NLP problem which will be discussed later.

A solution point whose objective value provides a valid lower bound for the optimum of problem (1), but not necessarily satisfying all constraints, is referred to as a dual solution. Like the standard OA algorithm, SCOT obtains dual solutions by solving relaxed problems that approximate the nonlinear constraints with polyhedral outer approximations. Depending on the type of outer approximations, the dual problem can be MILP, MIQP, or Mixed Integer Quadratically Constrained Linear (Quadratic) (MIQCL(Q)P) problems. Moreover, the dual bound is the best possible objective value of the dual problem. The primal and dual sub-problems are then iteratively solved by proper MIP algorithms. The MIP algorithms are distinguished depending on how the sub-problems are constructed, solved, and coordinated. Regardless of the solution algorithms adopted by SCOT, the dual and primal sub-problems are two primary components of the algorithms.

The main problem reformulation that SCOT attempts to solve by default is problem (6), which enforces both big-M and SOS-1 constraints and epigraph-reformulation. The separability of nonlinear functions in problem (6) allows SCOT to employ an alternative formulation, the so-called *lifted formulation* [22]. In this context, we use the lifted formulation for each nonlinear function, f_i , which results in tighter outer approximations when approximating nonlinear functions [19, 22, 37].

According to the idea of lifted formulation and also following the procedure presented in [22], at each iteration q , we construct the dual problem as the following MIP problem,

$$\min_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N \\ \mathbf{y}_1, \dots, \mathbf{y}_K \\ \gamma_1, \dots, \gamma_N}} \sum_{i=1}^N \gamma_i \quad (7a)$$

$$\text{subject to } \tilde{f}_i(\mathbf{x}_i^q) - \gamma_i \leq 0, \quad \forall \mathbf{x}_i^q \in \mathcal{X}_i^q, \quad \forall i = 1, \dots, N, \quad (7b)$$

$$\mathbf{A}\mathbf{x}_i \leq \mathbf{b}, \quad \forall i = 1, \dots, N, \quad (7c)$$

$$\mathbf{x}_i = \mathbf{y}_j, \quad \forall i \in \mathcal{E}_j, \quad \mathcal{E}_j \in \mathcal{E} \quad (7d)$$

$$-M_j \delta_{jk} \leq y_{jk} \leq M_j \delta_{jk}, \quad \forall j = 1, \dots, K, \quad \forall k = 1, \dots, n \quad (7e)$$

$$(y_{jk}, 1 - \delta_{jk}) : \text{SOS-1}, \quad \forall j = 1, \dots, K, \quad \forall k = 1, \dots, n \quad (7f)$$

$$\sum_{k=1}^n \delta_{jk} \leq \kappa, \quad \forall j = 1, \dots, K \quad (7g)$$

$$\delta_{jk} \in \{0, 1\}, \quad \forall j = 1, \dots, K, \quad \forall k = 1, \dots, n \quad (7h)$$

where $\gamma_i \in \mathcal{R}$, $i = \{1, \dots, N\}$, are new auxiliary decision variables, \mathbf{x}_i^q is a feasible point that satisfies linear and consensus constraints, obtained at iteration q , and $\tilde{f}_i(\mathbf{x}_i^q)$ is the outer approximation of $f_i(\mathbf{x}_i)$ around \mathbf{x}_i^q . Finally, \mathcal{X}_i^q is a finite set consisting of local feasible points defined as,

$$\mathcal{X}_i^q = \left\{ \mathbf{x}_i^\ell : \mathbf{A}\mathbf{x}_i^\ell \leq \mathbf{b}, \quad \forall \ell \in \{1, \dots, q\} \right\},$$

Set \mathcal{X}_i^q consists of the feasible points up to the current iteration q . The dual problem (7) is a relaxation of problem (6) since approximations of nonlinear constraints are used and its objective value is a lower bound of (6). At each iteration of SCOT algorithms, q , a new outer approximation is generated by each node of the network and

cooperatively added to problem (7) producing a tighter representation of the nonlinear constraints. The quality of outer approximations generated by $\tilde{f}_i(\mathbf{x}_i^q)$ directly impacts the convergence of the algorithms. Hence, SCOT provides first and second-order outer approximations and an event-triggered scheme that controls the effectiveness of the approximations. According to first-order Taylor series and convexity of $f_i(\mathbf{x}_i)$ functions, we can express constraints (7b) as,

$$f_i(\mathbf{x}_i^q) + \nabla f_i(\mathbf{x}_i^q)^T (\mathbf{x}_i - \mathbf{x}_i^q) - \gamma_i \leq 0, \quad (8)$$

which are linear inequalities. In case the nonlinear functions, $f_i(\mathbf{x}_i)$, are *strongly* convex functions, SCOT replaces constraints (7b) with the following quadratic inequalities,

$$f_i(\mathbf{x}_i^q) + \nabla f_i(\mathbf{x}_i^q)^T (\mathbf{x}_i - \mathbf{x}_i^q) + \frac{m_i^q}{2} \|\mathbf{x}_i - \mathbf{x}_i^q\|_2^2 - \gamma_i \leq 0 \quad (9)$$

where $m_i^q > 0$ is a constant such that $\nabla^2 f_i(\mathbf{x}_i) \succeq m_i^q I$. With $m_i^q > 0$, it is clear that the cut given by (9) is stronger than the cut given by (8). However, the quadratic cuts (9) tend to result in more challenging sub problems in BnB. Therefore, there can still be a computational advantage of the linear cuts.

Remark 1. For general strongly convex functions, m_i^q is not obtained easily. However, in some practical problems found in statistical learning and control, the computation of m_i^q is feasible. For example, the objective function in sparse Model Predictive Control (s-MPC) problems (which is a subclass of the SCO problem) is typically a convex quadratic function. For convex quadratic functions, m_i^q is the smallest Eigenvalue of the Hessian matrix. In machine learning problems the objective function usually consists of a convex function and a strongly convex *regularization* term. In this case, m_i^q can be computed from the regularization term.

A crucial step to forming the outer approximations is the computation of the approximation points x_i^q for which various strategies and methods exist. One of the well-known methods to obtain x_i^q is fixing the local binary decision variables of problem (6), $\delta_{jk} = \delta_{jk}^q$, and solving the resulting nonlinear optimization problem. In this case, a distributed convex NLP problem is solved, and its optimal solution provides a primal solution to the original MINLP problem. The problem of solving (6) for fixed binary variables is the primal problem of SCOT, which, at each iteration q , is defined as,

$$\min_{\substack{\gamma \\ \mathbf{x}_1, \dots, \mathbf{x}_N \\ \mathbf{y}_1, \dots, \mathbf{y}_K}} \sum_{i=1}^N \gamma_i \quad (10a)$$

$$\text{subject to } f_i(\mathbf{x}_i) - \gamma_i \leq 0 \quad (10b)$$

$$\mathbf{A}\mathbf{x}_i \leq \mathbf{b}, \forall i = 1, \dots, N, \quad (10c)$$

$$\mathbf{x}_i = \mathbf{y}_j, \forall i \in \mathcal{E}_j, \mathcal{E}_j \in \mathcal{E} \quad (10d)$$

$$-M_j \delta_{jk}^q \leq y_{jk} \leq M_j \delta_{jk}^q, \forall j = 1, \dots, K, \forall k = 1, \dots, n. \quad (10e)$$

Problem (10) is a distributed convex NLP problem and its solution has the advantage of generating linearizations about points closer to the feasible region. Therefore, primal solutions and primal bounds are obtained by iteratively solving problem (10).

In the case that the primal problem is a centralized NLP, a feasible point that satisfies all linear and nonlinear constraints is considered to be the primal solution candidate. However, when the primal problem has to be solved distributedly, as in problem (10), some numerical considerations have to be taken into account. Particularly, in this case, in addition to all linear and nonlinear constraints, the consensus constraints (10d) have to be satisfied which is more challenging to deal with since all computational nodes have to agree on a consensus solution. In case the primal solution does not satisfy the consensus constraints within an acceptable numerical tolerance, poor outer approximations are generated and added to the dual problem. Therefore a larger number of iterations are required by the distributed NLP solver, especially when a large computational network is considered. Another challenge in solving the primal problem distributedly is the communication burden between the nodes of the network and the LFCs.

4. Distributed Hybrid Outer Approximation

In this section, we propose the Distributed Hybrid Outer Approximation (DiHOA) algorithm to solve problem (6). In essence, DiHOA is developed based on the LP/NLP-based BnB algorithm proposed in [34]. The LP/NLP-based BnB algorithm is an implementation of the standard OA algorithm, where only a single BnB tree is built and outer approximations are added dynamically to the MIP master problem. Since only one BnB tree is constructed during the solution procedure, the LP/NLP-based BnB method is also called the single-tree OA algorithm. Similarly, the standard implementation of the OA algorithm where a BnB tree is constructed at each iteration of the algorithm is called multiple-tree OA.

Before discussing the DiHOA algorithm, we briefly present Distributed Primal Outer Approximation (DiPOA) which is proposed in [33]. DiPOA is the baseline algorithm developed to solve problem (6) which, in essence, extends the standard OA algorithm so that the main computational parts related to the NLP sub-problems are handled in a distributed fashion. From the numerical point of view, DiPOA is a hierarchical algorithm which consists of three main computational layers, namely, *primal*, *cutting-plane manager*, and *master* levels.

The primal level deals with the nonlinear optimization part, particularly problem (10), consisting of two sub-levels responsible for a specific computational task. The first sub-level aims to simultaneously solve multiple local nonlinear optimization problems in a fully decentralized fashion. The local nodes then *synchronously* communicate to the second sub-level, which is responsible for another phase of computations. The second sub-level is usually responsible for aggregating the solution of local NLP problems connected to each LFC by a series of unconstrained NLP problems. The solution of this level is then sent to the first sub-level until a consensus is reached.

The main purpose of the cutting-plane manager level is to generate and manage the OA linearizations, which are obtained around the generated feasible points provided by the primal level. This level is also responsible for adding the linearizations into the master's problem (7).

Finally, the master's level corresponds to solving problem (7) based on the cuts generated by the cutting-plane manager. The master MIP problem approximates the nonlinear functions of problem (6). As the number of cuts increases, this approximation improves until a good piece-wise outer approximator is achieved. The binary solution of the master level is then sent to the primal level, which, together with the nonlinear

optimization problems, solves another set of NLP problems.

Generally speaking, when solving problem (2) with DiPOA, most of the total solution time is usually spent on solving the MIP master’s problem. In such a scenario, the MIP problems are similar in consecutive iterations since they only differ by a few linear constraints. In particular, at iteration q of DiPOA, a new feasible point is provided by the primal, around which a new linear approximation constraint is generated and added to the master’s problem. In the next iteration, $q + 1$, the master’s problem is reconstructed and solved from scratch.

Based on what we proposed in [33] and according to [34], we develop the DiHOA algorithm to avoid constructing many similar MIP BnB trees. The main idea of DiHOA is to iteratively build a single branch-and-bound tree whereby the primal problem is solved distributedly, while dynamically updating the dual problem (7) without reconstructing the branch-and-bound tree. However, the single-tree OA algorithm may lead to a large BnB and weaker approximations. To avoid that, DiHOA introduces several second-order outer approximations of the nonlinear functions to the root of the BnB tree through an event-triggered scheme, which leads to a tighter problem representation and a BnB tree with a fewer number of nodes. This procedure constructs the first MIP dual problem and initiates the BnB algorithm. During the BnB search, as soon as a new integer-feasible solution is found, the primal problem (10) is distributedly solved to determine whether a *lazy constraint* removing this integer-feasible point should be generated. By lazy constraints, we mean cutting planes that are lazily added to the MIP model whenever an integer feasible solution is found. At this point, the RHADMM algorithm is applied, and the new primal information is distributed to the computational nodes of the network. Then the generated lazy constraint is added to the current node, and all open nodes of the BnB tree and the search continues. Therefore, it is not required to reconstruct the BnB tree as in multiple-tree algorithms and the same BnB tree can be used after adding new linearizations as lazy constraints. Finally, the algorithm is terminated until the MIP integer relaxation results in a feasible integer solution to the MINLP problem (6).

A detailed description of the DiHOA algorithm is summarized in Algorithm 1. It can be observed in Algorithm 1 that DiHOA consists of three primary computational phases, namely, `Initialization`, `MultipleTreeSearch`, and `SingleTreeSearch` steps. After the algorithm is successfully initialized, DiHOA starts a multiple-tree strategy with second-order outer approximations and continues the computations until either the solution is found or poor lower bound improvement is achieved. In the former case, the algorithm is terminated and the optimal solution is returned. In the latter case, however, DiHOA accumulates all the outer approximations obtained until iteration q in the root of the latest BnB tree and, then, it starts a single-tree search strategy whereby approximations are added dynamically. As the name of the algorithm suggests, DiHOA is a *hybrid* algorithm that combines both single-tree and multiple-tree strategies by introducing an event-triggered scheme that determines the switching iteration, q_{switch} , at which a single-tree search strategy is started. In the following, we describe each computational step in detail.

4.1. Initialization Step

The initialization step is started by solving the integer relaxation of problem (6) to construct the dual problem (7) in the first iteration of the algorithm. The integer relaxation of problem (6) is a consensus optimization problem that is solved using the

RHADMM algorithm. In case the relaxation obtains a feasible solution with respect to problem (6), we terminate DiHOA with the optimal solution. Otherwise, we generate N first-order outer approximations by using the local information available in each computational node of the network and construct the first MIP dual problem according to (7).

4.2. Multiple Tree Search Step

The multiple-tree search step of the DiHOA algorithm is a crucial step that directly affects the DiHOA performance. In the cases that q_{switch} is a large number, a pure multiple-tree algorithm with second-order outer approximations is obtained. Otherwise, the resulting algorithm becomes the single-tree OA. Therefore, q_{switch} should be determined in such a way that maximizes the DiHOA performance. To do so, we introduce an event-triggered scheme that selects q_{switch} based on the difference between two consecutive lower bounds (*i.e.*, $lb^q - lb^{q-1}$). In particular, during the solution procedure, DiHOA checks if the generated lower bounds by the dual problem start to flatten out within a given tolerance $\epsilon > 0$ and triggers a switching event, E_s , if $lb^q - lb^{q-1} \leq \epsilon$. As soon as E_s is triggered, DiHOA switches to the single-tree strategy by performing the BnB algorithm on the latest dual problem, which was obtained during the multiple-tree strategy.

4.3. Single Tree Search Step

The single-tree search step is activated for $q > q_{switch}$ after the E_s event is triggered. In this step, DiHOA accumulates all the outer approximations obtained during the `MultipleTreeSearch` phase in the root of the latest BnB tree and then starts the single-tree BnB procedure. The BnB search is initialized by solving an integer relaxation of the dual problem (7). In each node of the BnB tree, a QCLP relaxation is solved and the search is stopped once an integer solution is obtained in one of the nodes. The integer solution is then used to solve the primal problem (10) with integer variables fixed. The primal solution provides a valid upper bound and new approximations can be generated. The new approximations are then added to all open nodes in the BnB tree and the QCLP relaxation is resolved for the node which resulted in the integer combination. The BB procedure continues from the existing search tree with the improved polyhedral outer approximation.

As in the standard BnB, nodes can be pruned off in case the optimum of the QCLP relaxation exceeds the upper bound. However, the search cannot be stopped once an integer solution is obtained at a node, which must continue until the QCLP relaxation results in a feasible integer solution for problem (6) or until the node can be pruned off. Finally, since all variables in problem (1) are bounded, then the assumptions A1–A3 in [14] and assumptions in [21] are valid and Slater’s constraint qualification holds for problem (6) when the binary variables are fixed. Hence, no NLP subproblem will be infeasible and the convergence of Algorithm 1 is ensured.

5. SCOT: A Software Framework for Sparse Optimization

This section discusses the technical features of SCOT, including architecture and design, main components, basic syntax, and usage. SCOT is entirely written in C++17

Algorithm 1: DiHOA Algorithm

(1) Initialization

- 1.1 obtain a relaxed solution $\tilde{\mathbf{x}}_i^0$, $\tilde{\mathbf{y}}_j^0$, and $\tilde{\boldsymbol{\delta}}_j^0$, $i \in \{1, \dots, N\}$, $j \in \{1, \dots, K\}$ by solving an integer relaxation of the MINLP problem (6)
- 1.2 generate N outer approximations according to (8) and construct the dual problem (7)
- 1.3 store the generated outer approximations.
- 1.4 set $q = 1$, $ub^0 = +\infty$, $lb^0 = -\infty$

(2) MultipleTreeSearch**while** $lb^q - lb^{q-1} > \epsilon$ **do**

- 2.1 solve problem (7) and obtain a dual solution $\boldsymbol{\delta}_j^q$
- 2.2 $lb^q \leftarrow \sum_{i=1}^N \gamma_i$
- 2.3 solve problem (10) and obtain a primal solution \mathbf{x}_i^q , \mathbf{y}_j^q
- 2.4 $ub^q \leftarrow \min \left(\sum_{i=1}^N f_i(\mathbf{x}_i^q), ub^{q-1} \right)$
- 2.5 If $ub^q - lb^q \leq \tau$, return \mathbf{x}_i^q as an optimal solution \mathbf{x}_i^*
- 2.6 generate and store outer approximations according to (9) and update dual problem (7) by adding new outer approximations
- 2.7 $q \leftarrow q + 1$

end**(3) SingleTreeSearch**

- 3.1 start BnB search for the MIP problem obtained in the last step with all second-order outer approximations accumulated in the root

while $ub^q - lb^q > \epsilon$ **do**

- 3.2 if a new feasible integer solution, $\bar{\boldsymbol{\delta}}_j^q$, is found, then $\boldsymbol{\delta}_j^q \leftarrow \bar{\boldsymbol{\delta}}_j^q$ and update lb^q according to the current BnB tree
- 3.3 solve primal problem (10) and update the primal solution \mathbf{x}_i^q , \mathbf{y}_j^q
- 3.4 $ub^q \leftarrow \min \left(\sum_{i=1}^N f_i(\mathbf{x}_i^q), ub^{q-1} \right) \mathbf{x}_i^*$
- 3.5 generate outer approximations according to (8) and add them to the current and open nodes of BnB
- 3.6 resolve the integer relaxation problem for the node which resulted in the integer combination
- 3.7 $q \leftarrow q + 1$ and continue exploring BnB tree

end

programming language and at its core uses the *Message Passing Interface (MPI)* [16] to perform distributed computation operations, communication and message-passing between nodes of the network. To handle local NLP optimization problems, **SCOT** relies on various open source solvers, such as OSQP [35] and IPOPT [40]. **SCOT** also implements a truncated Newton’s method [29] to solve unconstrained optimization problems. Moreover, **SCOT** integrates commercial and open source MIP solvers such as **Gurobi** [18] and **CBC** to solve MIP dual problems. To solve the distributed NLP problem (10), **SCOT** implements the RHADMM algorithm proposed in [32] using the main MPI operations. Finally, both DiPOA and DiHOA algorithms are implemented within **SCOT** and are available through **SCOT** Python API, **ScotPy**, or **SCOT** Command Line Interface (CLI).

5.1. Architecture

A high-level overview of **SCOT** architecture, its main layers, and components are shown in Figure 1. As observed in the figure, the primary layers of the framework are **SCOTPY**, **SMCLI**, **SCOT Solver**, and **Computing Network** each of which consists of different tools and components to build and solve the optimization problem. The layered architecture of **SCOT** leads to a highly modular framework that can be easily extended by new features and algorithms. In the following, we describe each layer of **SCOT** and its components.

5.1.1. SCOTPY

SCOTPY is the main API written in **Python 3.8** programming language, which provides various modules and classes to define the optimization problem and solver settings. This layer consists of four components, namely, problem builder, validators and parsers, solver settings, and file generators. In principle, **SCOTPY** receives the problem input and algorithm settings from the application code layer and, after parsing and validating the input data and settings, writes the optimization problem and settings in **JSON** format with a specific naming convention. These steps are performed by validators and parsers, and file generator components. Moreover, according to the number of nodes given by the user, N objective functions with different problem data are created and stored on the file system as different **JSON** files.

Therefore, the resulting computation of this layer is to express the optimization problem and solver settings in various **JSON** files that can be read by the subsequent layers. Representing the optimization problem using files provides more flexibility since it decouples the optimization model from the optimization solver. Therefore, it is possible to write the optimization model in any language of choice and call the optimization solver from a different programming language or framework. Coupling the optimization model and the optimization solver through file formats is well-known in the optimization software industry and has been widely used for decades. We refer the interested reader to [23] for more details.

5.1.2. SMCLI

SCOT MPI Command Line Interface (**SMCLI**) is the main layer utilized to directly execute **SCOT Solver** according to different input and setting files. **SMCLI** can be used directly without **SCOTPY** interface, however the problem definition using **SCOTPY** is more appropriate. The validator and parsers component of **SMCLI** is responsible to

parse and validate the problem input and settings files, providing suitable data structures containing the optimization problem data for the `SCOT Solver`. Additionally, initializing computational nodes and various software libraries used in `SCOT Solver` are among the responsibilities of `SMCLI` layer.

5.1.3. *SCOT Solver*

At its core, the `SCOT` framework consists of `SCOT Solver` layer which is responsible to solve the optimization problem using a proper algorithm and settings. The main components of this layer are algorithms, engine, models, and utilities which are discussed in this section.

The `Engine` component consists of various modules and classes to distributedly solve the primal problem (10) and deliver the primal solution to the MINLP algorithms. Owing to its flexible and modular implementation, `Engine` can be easily extended by introducing user-defined and custom-distributed convex optimization algorithms and solvers. In principle, `Engine` requires the solution of local NLP problems for which multiple open-source and commercial solvers are available. At its core, `Engine` consists of a sub-module, `kernel` that provides a flexible interface to third-party solvers that can solve local optimization problems. The solvers supported by the `kernel` module are `OSQP`, `IPOPT`, and `Gurobi`. As a final note, the `Engine` component is called by all internal algorithms of `SCOT` and handles most of `MPI` communications and collective operations. More importantly, the quality of outer approximations depends on this component since it provides feasible points around which nonlinear functions are approximated.

One of the most critical components of `SCOT Solver` is `algorithm` component that implements Algorithms 1 and DiPOA. The `algorithms` component consists of two main modules, namely `dipoa`, and `dihoa`, which are responsible for implementing their corresponding algorithm. Because all the implemented algorithms must manage and monitor outer approximations, the `algorithm` component also includes a `managers` module. This module implements several classes to support the necessary data structures that generate and store both first- and second-order outer approximations. Moreover, the `managers` module implements the event-triggered schemes to improve the outer approximation's quality and switch from `MultipleTreeSearch` to `SingleTreeSearch` strategy. Among all classes, the `managers` module consists of two important classes, namely `CutStorage` and `CutGenerator`. `CutStorage` provides a simple way to validate and store the linear and quadratic outer approximations. The primary responsibility of the `CutGenerator` class is to generate necessary outer approximations from the information received from the `Engine` component. The `algorithm` module also provides a flexible functionality to interface third-party MIP solvers such as `Gurobi` and `Cplex` for solving the MIP dual problem (7).

The `model` component's primary responsibility is to generate a concrete internal representation of the optimization problem to be used by algorithms. This component consists of various classes to present different types of nonlinear objective functions and linear and sparsity constraints. By accessing the `model` component, the `algorithm` will be able to access the optimization problem data whenever necessary during the computations.

Finally, the `utilities` module implements classes and functions to provide commonly required functionalities, such as low-level parsing, measuring the CPU time, writing logs, constant parameters, exceptions, and file handling functions.

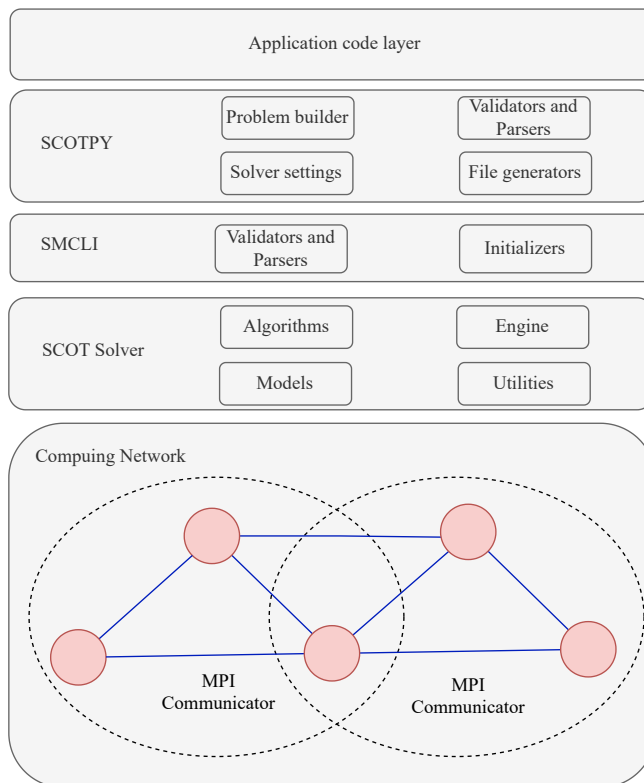


Figure 1. SCOT Architecture

5.1.4. Computing Network

This layer is responsible for presenting and managing the computational network using a graph data structure and Message Passing Interface (MPI) library. The **computing network** layer is in tight communication with the **SCOT Solver** layer since all distributed algorithms use MPI for performing distributed computations and inter-process communications.

5.2. Message Passing Interface (MPI)

This section provides an overview of MPI and its important operations. MPI is a message-passing library that supports parallel and distributed computations. For being independent of any programming language, MPI is arguably the most widely used platform for high-performance distributed computing nowadays. MPI is a software library for which various interfaces exist from various programming languages, including C/C++ and Python. MPI adopts the *Single Program, Multiple Data* (SPMD) programming paradigm to provide an efficient way to perform distributed computing. Using the SPMD paradigm, each node of the computing network runs the same program code, but it works with its own set of local variables and a separate subset of the data. To perform efficient distributed computing, MPI provides point-to-point and collective communications between the nodes of the computing network. Point-to-point communications refer to sending and receiving data between two different nodes, whereas collective communication is primarily performed among a set of nodes. In the following, we review some of the standard collective operations.

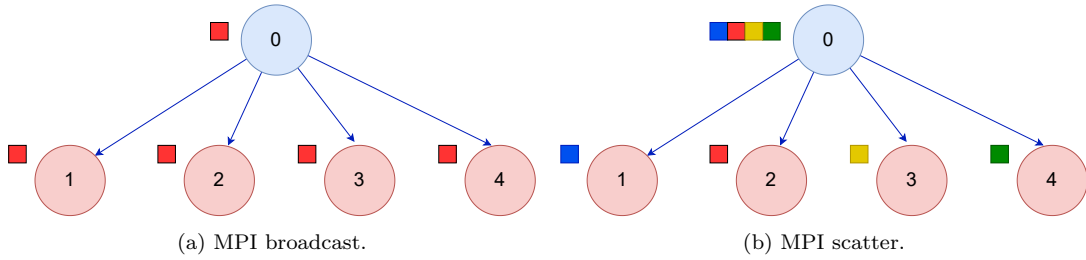


Figure 2. MPI broadcast and scatter communication patterns.

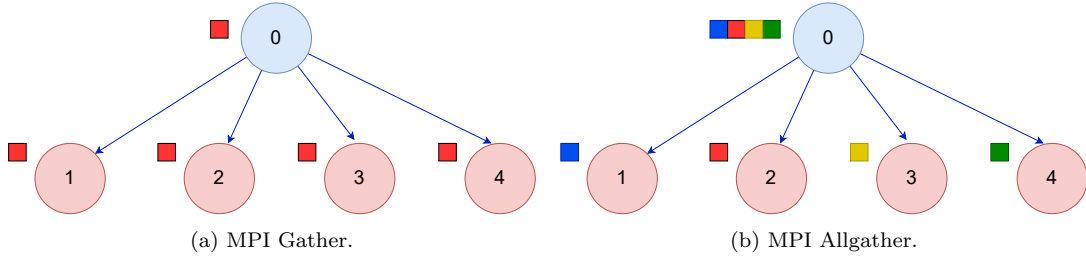


Figure 3. MPI Gather and Allgather communication patterns

5.2.1. MPI Broadcast

A *broadcast* is one of the basic collective communication strategies where one process sends the same data to all processes. Figure 2a illustrates the broadcast communication pattern.

5.2.2. MPI Scatter

A *scatter* is a collective routine that involves a designated root process sending data to all other processes. The main difference between MPI scatter and broadcast is that while the MPI broadcast sends the *same* piece of data to all other processes, the MPI scatter sends *chunks of an array* to different processes, as shown in Figure 2b.

5.2.3. MPI Gather and Allgather

In principle, the MPI *gather* is the inverse of the MPI scatter. Instead of distributing elements from one process to many processes, MPI gather takes elements from many processes and brings them together in one single process. The MPI *Allgather* collects all of the elements and then distributes them to all the processes. In the most basic scenario, MPI allgather is an MPI gather followed by an MPI broadcast. For example, Figure 3 illustrates the communication pattern in MPI gather and allgather.

5.2.4. MPI Reduce and Allreduce

MPI *reduce* involves reducing a set of elements into a small set of elements via a function. The MPI reduce takes an array of input elements from each process and returns an array of output elements to the root process. In a complementary style of MPI allgather to MPI gather, MPI allreduce will reduce the values and distribute the results to all processes. MPI reduce and allreduce communication patterns are

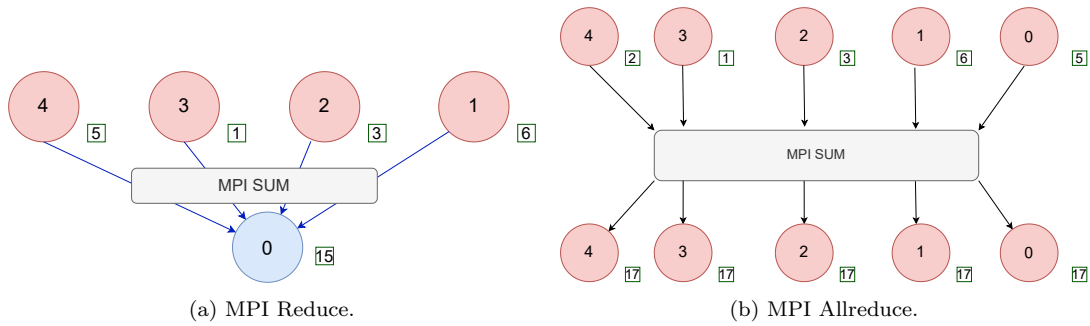


Figure 4. MPI Reduce and Allreduce Communication Patterns.

depicted in Figure 4.

5.3. Basic Syntax and Usage

In this section, we present an illustrative example to show how SCOT Python API, SCOTPY, is used to solve a distributed sparse logistic regression problem with random data. To do so, we first import the required classes from SCOTPY as the following code snippet shows,

```
from scotpy import (AlgorithmType,
                   ProblemType,
                   ScotModel,
                   ScotPy,
                   ScotSettings
                  )
```

Listing 1 Import statement of SCOT Python API

Here `ScotPy` is the main class that executes SCOT for a given problem and settings defined by `ScotModel` and `ScotSettings`, respectively. The `AlgorithmType` and `ProblemType` classes determine what problem class is solved and which algorithm will be used. In order to create the optimization problem and SCOT settings, the following code snippet can be used,

```
# Create a classification dataset with 1000 rows and 20 columns
dataset, res = make_classification(n_samples = 1000, n_features = 20)
scp = ScotModel(problem_name = "logistic_regression",
               rank = 0,
               kappa = 5,
               ptype = ProblemType.CLASSIFICATION)

# Set problem data with normalization
scp.set_data(dataset, res, normalized_data = True)

# Create corresponding JSON files that represent the optimization problem.
scp.create()

scot_settings = ScotSettings(
    relative_gap = 1e-5,
    time_limit = 100,
    verbose = True,
    algorithm = AlgorithmType.DIHOA
)
```

Listing 2 Problem definition and settings

where `make_classification` function, imported from Python `scikit-learn` library, is used to generate a random classification dataset. The `ScotModel` object is then created by a given problem name, MPI rank, number of nonzeros, and problem type. The solver settings can be defined by creating an object from `ScotSettings` class. We note that by executing MPI, the above code snippets are simultaneously executed by each node of the network. Hence, each node can use its own problem data. Finally, we solve the optimization problem by using the following code,

```
solver = ScotPy(problem, scot_settings)
status_code = solver.run()
```

Listing 3 SCOT Execution

where `ScotPy` class is responsible for creating a solver object for a given problem and settings. By executing the `run` method of `ScotPy`, MPI execution with N nodes is started.

6. Applications and Numerical Examples

In this section, we evaluate the SCOT performance by comparing SCOT to various state-of-the-art MINLP solvers equipped with single-tree, multiple-tree, and nonlinear BnB algorithms. We considered SHOT and BONMIN as decomposition-based and KNITRO as nonlinear BnB solvers. However, one should keep in mind that these are general-purpose solvers, and unlike SCOT they are not tailored for this specific problem structure. From the application point of view, we focus on sparse logistic and linear regression problems with data distributed over a network of computational nodes. The benchmark results are obtained by generating performance profiles according to various problem instances over different computational nodes generated by SCOTPY.

6.1. Implementation Details and Set-up

All the experiments were performed on a Linux machine with an Intel Core i5 2.50 GHz processor, with four physical cores and 16 GB of RAM. The SCOT solver is entirely implemented in C++ programming language and relies on MPI to carry out distributed computation and inter-process communication. The source code of the solver is available on <https://github.com/Alirezalm/scot>. To perform linear algebra operations required by the distributed NLP solver, SCOT uses Eigen 3.4 library. Moreover, the MIP solver employed in SCOT is GUROBI 9.5.2 with an academic license. As for the comparison with other MINLP solvers, GAMS 36.2 was selected as the optimization platform. It should also be noted that to achieve a meaningful comparison, GUROBI is selected as the primary MIP solver for all MINLP solvers considered in the benchmarks.

6.2. Experimental Set-up

This section presents the problem classes that compose the benchmarks and experiments. We consider two well-known machine-learning problems: classification and regression. For both problem classes, the dataset is assumed to be distributed over a computational network, and the sparsity of the solution matters. Sparse classification and regression problems are among the central problems in statistics and machine learning, whose solutions usually lead to more interpretable models. In this paper, we

choose logistic regression and linear regression as candidate models for classification and regression.

Sparse classification and regression problems are tightly connected to SCO problems as it is often desired to identify a critical subset of features contributing to the response. Furthermore, the sparse solution usually leads to more interpretable models and improves prediction accuracy by eliminating unnecessary features. Accordingly, we introduce Distributed Sparse Logistic Regression (DSLogR) and Distributed Sparse Linear Regression (DSLInR) problems. The DSLogR problem is defined as,

$$\min_{\boldsymbol{\theta}_i, \mathbf{y}_j, \delta_j} \sum_{i=1}^N \sum_{\ell=0}^p \log \left[1 + e^{-(\boldsymbol{\theta}_i^T \mathbf{X}_{i,\ell}) \Gamma_{i,\ell}} \right] + \frac{\lambda}{2} \|\boldsymbol{\theta}_i\|_2^2 \quad (11a)$$

$$\boldsymbol{\theta}_i = \mathbf{y}_j, \forall i \in \mathcal{E}_j, \mathcal{E}_j \in \mathcal{E}, \forall j = 1, \dots, K \quad (11b)$$

$$-M_j \delta_{jk} \leq y_{jk} \leq M_j \delta_{jk}, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (11c)$$

$$\sum_{k=1}^n \delta_{jk} \leq \kappa, \forall j = 1, \dots, K \quad (11d)$$

$$y_{jk}(1 - \delta_{jk}^q) = 0, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (11e)$$

$$\delta_{jk} \in \{0, 1\}, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (11f)$$

where $\mathbf{X}_i \in \mathcal{R}^{p \times n}$ and $\boldsymbol{\Gamma}_i \in \mathcal{R}^p$ are the dataset and response vector of the i -th node, $\mathbf{X}_{i,\ell}$ is the column representation of ℓ -th row of \mathbf{X}_i , $\Gamma_{i,\ell}$ is the ℓ -th element of $\boldsymbol{\Gamma}_i$, and $\lambda > 0$ is the regularization parameter. Clearly, problem (11) is a subclass of the SCO problem (6) which is the primary reformulation that SCOT uses. Similarly, the DSLInR problem is defined as,

$$\min_{\boldsymbol{\theta}_i, \mathbf{y}_j, \delta_j} \sum_{i=1}^N \|\mathbf{X}_i \boldsymbol{\theta}_i - \mathbf{b}_i\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}_i\|_2^2 \quad (12a)$$

$$\boldsymbol{\theta}_i = \mathbf{y}_j, \forall i \in \mathcal{E}_j, \mathcal{E}_j \in \mathcal{E}, \forall j = 1, \dots, K \quad (12b)$$

$$-M_j \delta_{jk} \leq y_{jk} \leq M_j \delta_{jk}, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (12c)$$

$$\sum_{k=1}^n \delta_{jk} \leq \kappa, \forall j = 1, \dots, K \quad (12d)$$

$$y_{jk}(1 - \delta_{jk}^q) = 0, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (12e)$$

$$\delta_{jk} \in \{0, 1\}, \forall j = 1, \dots, K, \forall k = 1, \dots, n \quad (12f)$$

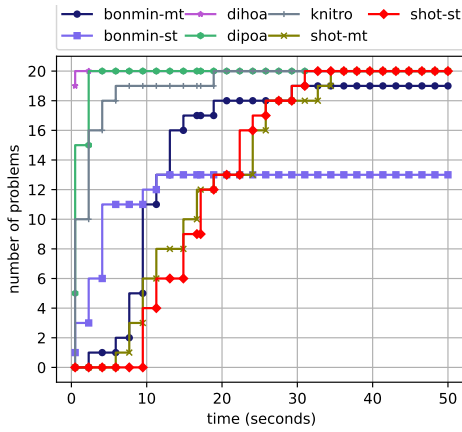
We generate N random local datasets for both DSLogR and DSLInR problems with zero mean and unit ℓ_2 norm for each column. We perform the numerical benchmarks based on different solver settings. Default settings were adopted for Knitro as an NLP BnB solver. The chosen settings for SCOT, SHOT, and Bonmin are reported in Table 2.

6.3. Benchmark Results

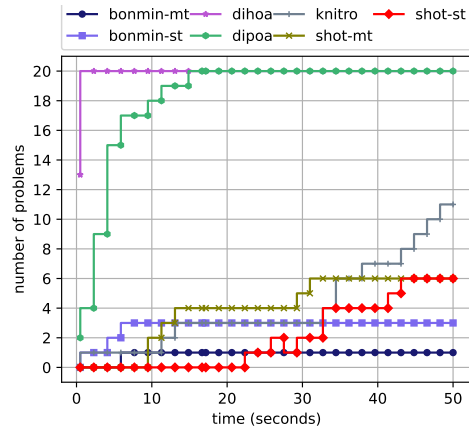
This section presents performance profiles comparing the SCOT performance with some MINLP solvers with different settings. We considered seven benchmark scenarios containing 20 different problem instances with different properties and settings. In each scenario, we generate 15 DSLogR and 5 DSLInR random problem instances with a

Table 2. MINLP Solver Settings

solver	Settings			
	algorithm	name	mip solver	nlp solver
SCOT	dipoa	scot-mt	gurobi	rhadmm
SCOT	dihoa	scot-st	gurobi	rhadmm
SHOT	ESH multiple-tree	shot-mt	gurobi	ipopt
SHOT	ESH single-tree	shot-st	gurobi	ipopt
Bonmin	B-OA	bonmin-mt	gurobi	ipopt
Bonmin	B-QG	bonmin-st	gurobi	ipopt

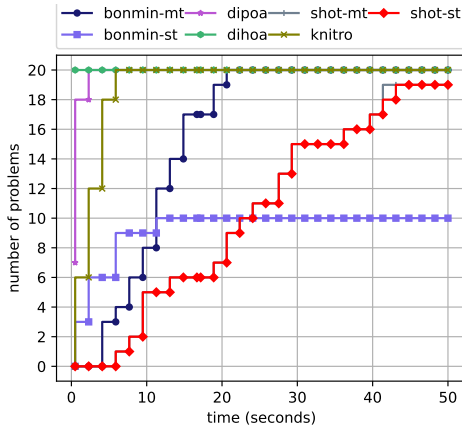


(a) 90% sparsity.

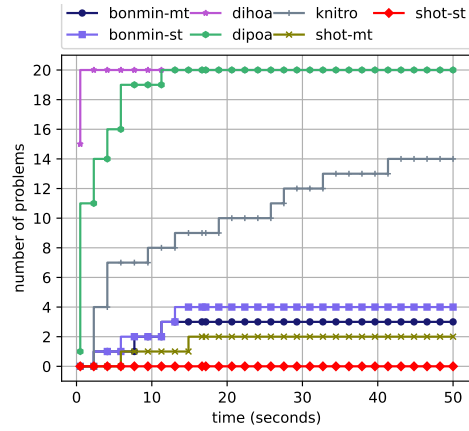


(b) 80% sparsity.

Figure 5. Benchmark results for scenario 1.

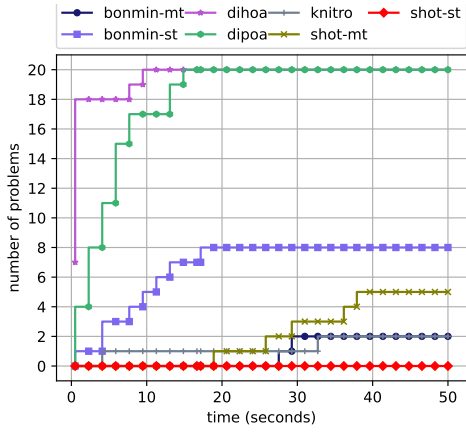


(a) 90% sparsity.

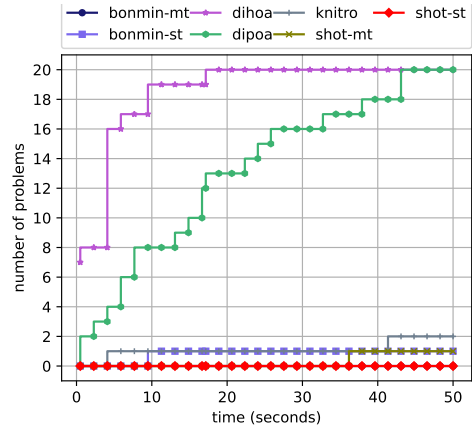


(b) 80% sparsity.

Figure 6. Benchmark results for scenario 2.

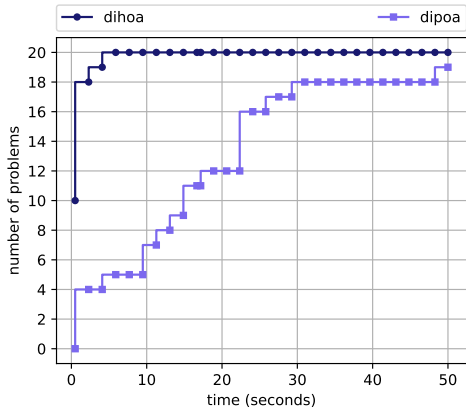


(a) 90% sparsity.

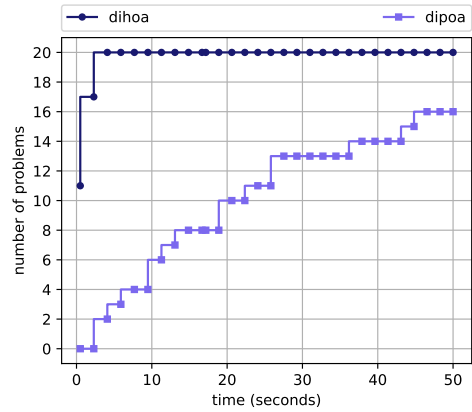


(b) 80% sparsity.

Figure 7. Benchmark results for scenario 3.

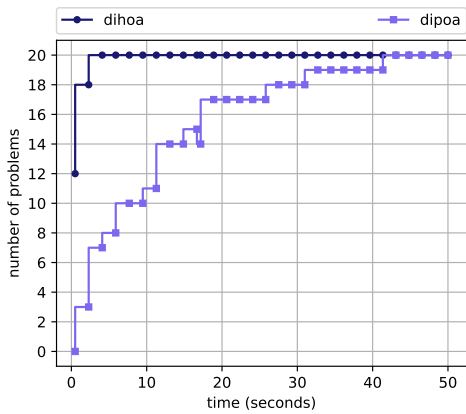


(a) 90% sparsity.

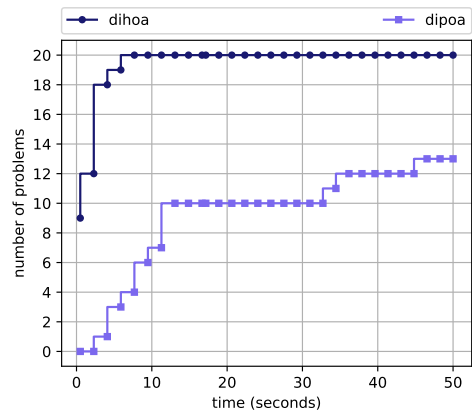


(b) 80% sparsity.

Figure 8. Benchmark results for scenario 4.

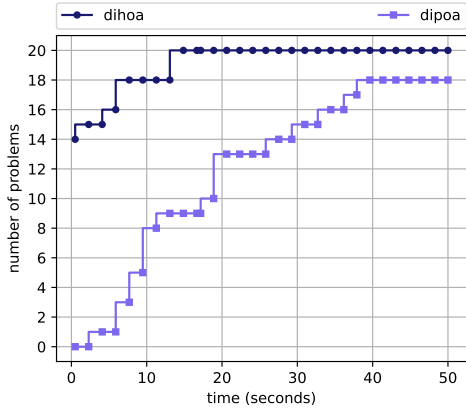


(a) 90% sparsity.

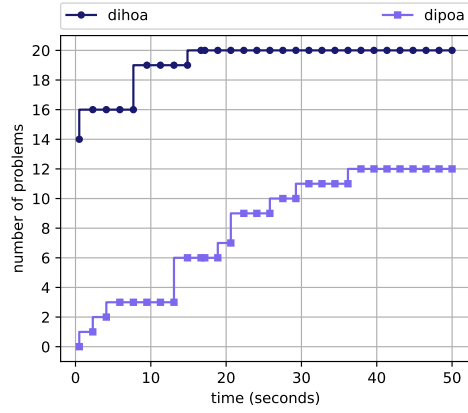


(b) 80% sparsity.

Figure 9. Benchmark results for scenario 5.

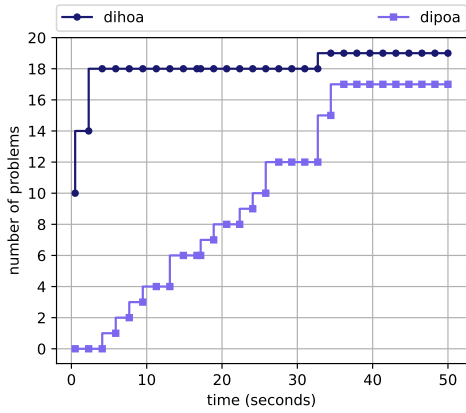


(a) 90% sparsity.

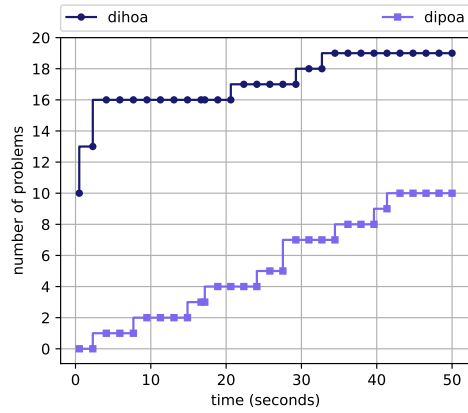


(b) 80% sparsity.

Figure 10. Benchmark results for scenario 6.



(a) 90% sparsity.



(b) 80% sparsity.

Figure 11. Benchmark results for scenario 7.

Table 3. Benchmark Settings for Performance Profiles

scenario	Scenario settings						
	n_{\min}	n_{\max}	p_{\min}	p_{\max}	N	n_p	p_{tot}
1	20	30	1000	2500	2	20	5000
2	20	30	1000	5000	2	20	10000
3	25	50	1000	5000	2	20	10000
4	25	100	1000	10000	4	20	40000
5	25	100	1000	20000	4	20	80000
6	25	100	1000	50000	4	20	200000
7	25	200	1000	50000	6	20	300000

different number of features and sample points within a given range. Therefore the benchmark set consists of a total of 140 problem instances. Moreover, each algorithm appearing in Table 2 is applied to solve all problem instances with 30 different maximum execution times limits, starting from 0.5 to 50 seconds. Therefore, the total number of algorithm runs for each scenario is 600, leading to 4200 algorithm executions for all scenarios. Table 3 represents settings for each benchmark scenario where n_{\min} and n_{\max} are the minimum and the maximum number of features, p_{\min} and p_{\max} are the minimum and the maximum number of data points for each computational node, n_p is the total number of problems, and p_{tot} is the total number of data points considering all computational nodes.

Figures 5-11 compare the solvers **SCOT**, **Bonmin**, **SHOT**, and **Knitro** with both single-tree and multiple-tree algorithms. The comparison results for each scenario are shown as performance profiles consisting of two different sparsity levels of the solution.

The benchmark results of the first three scenarios are depicted in Figures 5-7 where *small to medium size* problem instances are considered. In both sparsity levels, the DiPOA and DiHOA algorithms show better performance compared to other MINLP solvers. It can be observed in Figure 7 that for larger problem instances the performance gap between **SCOT** and other MINLP solvers is increased.

Figures 8-11 depict the performance profiles for the scenarios 4-7 in which *medium to large* problem instances are taken into account. In these scenarios, all MINLP solvers failed to solve the problem with the considered maximum execution time limit. Therefore, we only provide performance profiles of DiPOA and DiHOA algorithms. According to the performance profiles, the DiHOA algorithm is more efficient compared to DiPOA as the performance gap between them is increasing for large problem instances. For example, in scenario 7 with 80% of sparsity, DiHOA was able to solve 19 problem instances within the given maximum execution time limit, whereas DiPOA only solved 10 instances.

6.4. Discussion

In this section, we compared the performance of **SCOT** using both DiHOA and DiPOA algorithms with some state-of-the-art MINLP solvers. According to the numerical experiments and performance profiles, the DiHOA algorithm achieves better performance and efficiency in all problem instances. However, one should keep in mind that **SCOT** is tailored for SCO problems which is not the case for the general purpose solvers. The results also showed that for large and distributed problems a distributed solver can provide an efficient and reliable solution.

7. Conclusion

In this work, we introduced SCOT solver and DiHOA algorithm that SCOT implements to solve the DSCO problem (1). The DiHOA implements an outer approximation algorithm to solve an MINLP equivalent to the SCO problem, combining single- and multiple-tree outer approximation with a decentralized algorithm for solving convex nonlinear subproblems, which generates primal solutions and cuts for the master dual problem. The numerical benchmark results and comparison to state-of-the-art MINLP solvers indicated that SCOT equipped with DiHOA can be efficiently used for solving sparse convex optimization problems in different application domains. The performance and efficiency of the solver were achieved by augmenting SCOT with modern convex and mixed-integer optimization techniques. Continuing work aims at enhancing SCOT effectiveness by means of the development of decentralized mixed integer algorithms and new heuristic techniques to obtain outer approximations.

Disclosure statement

The authors report there are no competing interests to declare.

Funding

This work was funded in part by Fundação de Amparo à Pesquisa e Inovação do Estado de Santa Catarina (FAPESC) under grant 2021TR2265 and Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES, Brazil) under the project PrInt CAPES-UFSC 698503P1. Financial support from Digital Futures at KTH is also gratefully acknowledged.

References

- [1] R.P. Aguilera, G. Urrutia, R.A. Delgado, D. Dolz, and J.C. Agüero, *Quadratic model predictive control including input cardinality constraints*, IEEE Transactions on Automatic Control 62 (2017), pp. 3068–3075.
- [2] D. Bertsimas and R. Cory-Wright, *A scalable algorithm for sparse portfolio selection*, INFORMS Journal on Computing (2022).
- [3] D. Bertsimas, J. Dunn, L. Kapelevich, and R. Zhang, *Sparse regression over clusters: SparClur*, Optimization Letters (2021), pp. 1–16.
- [4] D. Bertsimas and A. King, *Logistic regression: From art to science*, Statistical Science 32 (2017), pp. 367–384.
- [5] D. Bertsimas, A. King, and R. Mazumder, *Best subset selection via a modern optimization lens*, The annals of statistics 44 (2016), pp. 813–852.
- [6] D. Bertsimas and N. Mundru, *Sparse convex regression*, INFORMS Journal on Computing 33 (2021), pp. 262–279.
- [7] D. Bertsimas, J. Pauphilet, and B. Van Parys, *Sparse classification: a scalable discrete optimization perspective*, Machine Learning 110 (2021), pp. 3177–3209.
- [8] D. Bertsimas and R. Weismantel, *Optimization over integers*, Dynamic Ideas Belmont, MA, 2005.
- [9] D. Bienstock, *Computational study of a family of mixed-integer quadratic programming problems*, Mathematical Programming 74 (1996), pp. 121–140.

- [10] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *et al.*, *Distributed optimization and statistical learning via the alternating direction method of multipliers*, Foundations and Trends® in Machine Learning 3 (2011), pp. 1–122.
- [11] S.S. Chen, D.L. Donoho, and M.A. Saunders, *Atomic decomposition by basis pursuit*, SIAM review 43 (2001), pp. 129–159.
- [12] M.A. Duran and I.E. Grossmann, *An outer-approximation algorithm for a class of mixed-integer nonlinear programs*, Mathematical Programming 36 (1986), pp. 307–339.
- [13] Y.C. Eldar and G. Kutyniok, *Compressed sensing: theory and applications*, Cambridge university press, 2012.
- [14] R. Fletcher and S. Ley, *Solving Mixed Integer Nonlinear Programs by Outer Approximation 1 Introduction*, October 66 (1996), pp. 327–349.
- [15] S. Foucart and H. Rauhut, *An invitation to compressive sensing*, in *A Mathematical Introduction to Compressive Sensing*, Springer, 2013, pp. 1–39.
- [16] W. Gropp, E. Lusk, and A. Skjellum, *Using MPI: Portable Parallel Programming with the Message-Passing Interface*, Vol. 1, MIT Press, 1999.
- [17] I. Grossmann, *Review of Nonlinear Mixed-Integer and Disjunctive Programming Techniques*, Optimization and Engineering 3 (2002), pp. 227–252.
- [18] L. Gurobi Optimization, *Gurobi optimizer reference manual* (2020). Available at <http://www.gurobi.com>.
- [19] H. Hijazi, P. Bonami, and A. Ouorou, *An outer-inner approximation for separable mixed-integer nonlinear programs*, INFORMS Journal on Computing 26 (2014), pp. 31–44.
- [20] Y. Kim, J. Kim, and Y. Kim, *Blockwise sparse regression*, Statistica Sinica (2006), pp. 375–390.
- [21] J. Kronqvist, D.E. Bernal, and I.E. Grossmann, *Using regularization and second order information in outer approximation for convex MINLP*, Mathematical Programming 180 (2020), pp. 285–310.
- [22] J. Kronqvist, A. Lundell, and T. Westerlund, *Reformulations for utilizing separability when solving convex MINLP problems*, Journal of Global Optimization 71 (2018), pp. 571–592.
- [23] B. Legat, O. Dowson, J.D. Garcia, and M. Lubin, *MathOptInterface: a data structure for mathematical optimization problems*, INFORMS Journal on Computing 34 (2021), pp. 672–689.
- [24] F.L. Lewis, *Wireless sensor networks*, Smart environments: technologies, protocols, and applications (2004), pp. 11–46.
- [25] A. Lundell and J. Kronqvist, *Integration of polyhedral outer approximation algorithms with MIP solvers through callbacks and lazy constraints*, AIP Conference Proceedings 2070 (2019).
- [26] S. Ma, X. Song, and J. Huang, *Supervised group lasso with applications to microarray data analysis*, BMC bioinformatics 8 (2007), pp. 1–17.
- [27] B.K. Natarajan, *Sparse approximate solutions to linear systems*, SIAM journal on computing 24 (1995), pp. 227–234.
- [28] A. Nedić and J. Liu, *Distributed optimization for control*, Annual Review of Control, Robotics, and Autonomous Systems 1 (2018), pp. 77–103.
- [29] J. Nocedal and S. Wright, *Numerical optimization*, Springer Science & Business Media, 2006.
- [30] I. Notarnicola, Y. Sun, G. Scutari, and G. Notarstefano, *Distributed big-data optimization via block-iterative convexification and averaging*, in *IEEE 56th Annual Conference on Decision and Control (CDC)*, Dec. 2017, pp. 2281–2288.
- [31] G. Notarstefano, I. Notarnicola, A. Camisa, *et al.*, *Distributed optimization for smart cyber-physical networks*, Foundations and Trends® in Systems and Control 7 (2019), pp. 253–383.
- [32] A. Olama, N. Bastianello, P.R. Mendes, and E. Camponogara, *Relaxed hybrid consensus ADMM for distributed convex optimisation with coupling constraints*, IET Control Theory & Applications 13 (2019), pp. 2828–2837.

- [33] A. Olama, E. Camponogara, and P. Mendes, *Distributed primal outer approximation algorithm for sparse convex programming with separable structures*, arXiv preprint arXiv:2210.06913 (2021).
- [34] I. Quesada and I.E. Grossmann, *An LP/NLP based branch and bound algorithm for convex MINLP optimization problems*, Computers & Chemical Engineering 16 (1992), pp. 937–947.
- [35] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, *OSQP: An operator splitting solver for quadratic programs*, Mathematical Programming Computation 12 (2020), pp. 637–672.
- [36] X. Sun, X. Zheng, and D. Li, *Recent advances in mathematical programming with semi-continuous variables and cardinality constraint*, Journal of the Operations Research Society of China 1 (2013), pp. 55–77.
- [37] M. Tawarmalani and N.V. Sahinidis, *A polyhedral branch-and-cut approach to global optimization*, Mathematical programming 103 (2005), pp. 225–249.
- [38] R. Tibshirani, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society: Series B (Methodological) 58 (1996), pp. 267–288.
- [39] A.M. Tillmann, D. Bienstock, A. Lodi, and A. Schwartz, *Cardinality minimization, constraints, and regularization: A survey*, arXiv preprint arXiv:2106.09606 (2021).
- [40] A. Wächter and L.T. Biegler, *On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming*, Mathematical Programming 106 (2006), pp. 25–57.