

Distributed Algorithms for Robust Convex Optimization via the Scenario Approach

Keyou You, *Member, IEEE*, Roberto Tempo, *Fellow, IEEE*, and Pei Xie

Abstract—This paper proposes distributed algorithms to solve robust convex optimization (RCO) when the constraints are affected by nonlinear uncertainty. We adopt a scenario approach by randomly sampling the uncertainty set. To facilitate the computational task, instead of using a single centralized processor to obtain a “global solution” of the scenario problem (SP), we resort to *multiple interconnected processors that are distributed among different nodes of a network to simultaneously solve the SP*. Then, we propose a primal-dual sub-gradient algorithm and a random projection algorithm to distributedly solve the SP over undirected and directed graphs, respectively. Both algorithms are given in an explicit recursive form with simple iterations, which are especially suited for processors with limited computational capability. We show that, if the underlying graph is strongly connected, each node asymptotically computes a common optimal solution to the SP with a convergence rate $O(1/(\sum_{t=1}^k \zeta^t))$ where $\{\zeta^t\}$ is a sequence of appropriately decreasing stepsizes. That is, the RCO is effectively solved in a distributed way. The relations with the existing literature on robust convex programs are thoroughly discussed and an example of robust system identification is included to validate the effectiveness of our distributed algorithms.

Index Terms—Robust convex optimization, uncertainty, scenario approach, primal-dual algorithm, random projection algorithm.

I. INTRODUCTION

A robust convex optimization (RCO) is a convex optimization problem where an infinite number of constraints are parameterized by uncertainties. This problem has found wide applications in control analysis and synthesis of complex systems, as well as in other areas of engineering [1], [2]. As the dependence of the constraints on the uncertainties may be nonlinear, RCO is generally not easily solvable. In fact, the study of RCO bears a vast body of literature, see e.g. [3]–[5] and references therein.

In this paper, we adopt a *scenario approach*, which was first introduced in [1], [6] to solve RCO. In particular, we randomly sample the uncertainty set and obtain a standard convex optimization called the scenario problem (SP). The guarantees of optimality are then given in a probabilistic sense and an explicit bound on the probability that the original constraints are violated is provided. The striking feature of this approach is that the sample complexity, which guarantees that

a solution to the SP is optimal with a given level of confidence, can be computed a priori. We also refer to [2], [7] for general properties and specific randomized algorithms to cope with uncertainty in systems and control.

To facilitate the computational task, instead of using a single processor to solve the SP, this paper proposes a distributed computing framework with many interconnected processors. The challenging problem is to distribute the computational task among the nodes of a network, each representing a single processor. The idea is to break a (possibly) large number of constraints of the SP into many small sets of *local* constraints that can be easily handled in each node. That is, each node computes some optimal solution of the SP with a low computational cost. Under local interactions between nodes, the SP is then collaboratively solved in every node via three key steps.

First, every node randomly samples the uncertainty set of RCO, with the sample size inversely proportional to the total number of nodes or being a priori determined by its computational capability. Although this idea has been adopted in [8], [9] to solve the SP, our approach is substantially different. In particular, after sampling, each node in [8] requires to completely solve a local SP at each iteration and exchange the set of active constraints with its neighbors. The process continues until a consensus on the set of active constraints is reached. Finally, every node solves its local SP under *all* active constraints of the SP. Clearly, the number of constraints in every local SP increases with the number of iterations. In some extreme cases, each constraint in the SP can be active, and every node eventually solves a local SP that has the same number of constraints as the SP. Thus, the computational cost in each node is not necessarily reduced. Moreover, each node cannot guarantee to obtain the same optimal solution to the SP. Since an active constraint may become inactive in any future iteration, identifying the active constraints cannot be recursively implemented, and this computation is very sensitive to numerical errors. On the contrary, each node in this paper only needs to handle a fixed number of local constraints and recursively run an explicit algorithm with very simple structure.

Second, the SP is reformulated as a distributed optimization problem with many decoupled small sets of local constraints and a coupled constraint, which is specially designed in conformity with the network structure. If the number of nodes is large, each node only needs to deal with a very small number of local constraints. The information is then distributed across the network via the coupled constraint, so that it can be locally handled. We recall that a similar technique has been already

This work was supported by the National Natural Science Foundation of China (61304038), Tsinghua University Initiative Scientific Research Program.

Keyou You and Pei Xie are with the Department of Automation and TNList, Tsinghua University, 100084, China (email: youky@tsinghua.edu.cn, xie-p13@mails.tsinghua.edu.cn).

Roberto Tempo is deceased and was with CNR-IEIIT, Politecnico di Torino, Torino, 10129, Italy.

adopted to solve distributed optimization problems, see e.g. [10], [11], which are only focused on convex optimization problems and no robustness issues are addressed. On the other hand, robust optimization has also attracted significant attention in many research areas [12], [13], but the proposed approaches are fully centralized. In this paper, we address both distributed and robust optimization problems simultaneously.

Third, each node of the network keeps updating a local copy of an optimal solution by individually handling its local constraints and interacting with its neighbors to address the coupled constraint. If the graph is strongly connected, every pair of nodes can indirectly access information from each other. To this purpose, we develop two recursive distributed algorithms for each node to interact with the neighbors to solve the SP by utilizing the constraint functions under undirected and directed graphs, respectively. For both algorithms, the computational cost per iteration only involves a few additions and multiplications of vectors, in addition to the computation of the sub-gradients of parameterized constraint functions. Thus, the computational cost is small in each node, and the approach is particularly useful for solving a large-size optimization problem with many solvers of reduced power.

For undirected graphs, where the information flow between the nodes is bidirectional, we solve the distributed optimization problem by using an augmented Lagrangian function with a quadratic penalty [14]. Following this approach, a distributed primal-dual sub-gradient algorithm is designed to find a saddle point. In this case, both the decoupled and coupled constraints are handled by introducing Lagrange multipliers, which provide a natural approach from the optimization viewpoint. For the coupled constraint, each node also needs to broadcast its estimate of an optimal solution to the SP, and the modified Lagrange multipliers to the neighbors, after which it recursively updates them by jointly using sub-gradients of local constraint functions. We show that each node finally converges to some common optimal solution to the SP. We remark that most of the existing work on distributed optimization [11], [15], [16] uses the Euclidean projection to handle local constraints. The projection is easy to perform only if the projection set has a special structure, which is generally not the case in the SP. From this perspective, our algorithm is more attractive to solve the SP problem in the context of distributed algorithms.

For directed graphs, the information flow between nodes is unidirectional and the primal-dual algorithm for undirected graphs cannot be used. To overcome this issue, we address the coupled constraint by adopting a consensus algorithm and design a novel two-stage recursive algorithm. At the first stage, we solve an unconstrained optimization problem which removes the decoupled local constraints in the reformulated distributed optimization and obtain an intermediate state vector in each node. We notice that, in the classical literature [16]–[19], the assumption on balanced graphs is often made. In our paper, this restrictive assumption is removed and this step is non-trivial, see e.g. [20], [21]. At the second stage, each node individually addresses its decoupled local constraints by adopting a generalization of Polyak random algorithm [22], which moves its intermediate state vector toward a randomly selected local constraint set. Combining these two stages, and

under some mild conditions, both consensus and feasibility of the iteration in each node are achieved almost surely. Although this distributed algorithm is completely different from the primal-dual sub-gradient algorithm previously described, both algorithms essentially converge at a speed $O(1/(\sum_{t=1}^k \zeta^t))$ where $\{\zeta^t\}$ is a sequence of appropriately decreasing stepsizes.

The rest of this paper is organized as follows. In Section II, we formulate RCO and include four motivating examples, after which the probabilistic approach to RCO is introduced. In Section III, we describe a distributed computing framework for the SP. In Section IV, a distributed algorithm is proposed via the primal-dual sub-gradient method for undirected graphs and show its convergence. In Section V, we design a distributed random projected algorithm over directed graphs to solve RCO. An example focused on robust system identification is included in Section VI. Some brief concluding remarks are drawn in Section VII.

A preliminary version of this work appeared in [23], which only addresses undirected graphs with a substantially different approach. This paper provides significant extensions to directed graphs using randomized algorithms, establish their convergence properties, include the complete proofs and provide new simulation results for robust system identification.

Notation: The sub-gradient of a vector function $y = [y_1, \dots, y_n]' \in \mathbb{R}^n$ whose components are convex functions with respect to an input vector $x \in \mathbb{R}^m$ is denoted by $\partial y = [\partial y_1, \dots, \partial y_n]' \subseteq \mathbb{R}^{n \times m}$. For two non-negative sequences $\{a^k\}$ and $\{b^k\}$, if there exists a positive constant c such that $a^k \leq c \cdot b^k$, we write $a^k = O(b^k)$. For two vectors $a = [a_1, \dots, a_n]'$ and $b = [b_1, \dots, b_n]'$, the notation $a \succeq b$ means that a_i is greater than b_i for any $i \in \{1, \dots, n\}$. A similar notation is used for \succ , \preceq and \prec . The symbol $\mathbf{1}$ denotes the vector with all entries equal to one. Given a pair of real matrices of suitable dimensions, \otimes indicates their Kronecker product. Finally, $f(\theta)_+ = \max\{0, f(\theta)\}$ is the positive part of f , $\text{Tr}(\cdot)$ is the trace of a matrix and $\|\cdot\|$ denotes Euclidean norm.

II. ROBUST CONVEX OPTIMIZATION AND SCENARIO APPROACH

A. Robust Convex Optimization

Consider a robust convex optimization (RCO) of the form

$$\min_{\theta \in \Theta} c'\theta \quad \text{subject to } f(\theta, q) \leq 0, \forall q \in \mathcal{Q}, \quad (1)$$

where $\Theta \subseteq \mathbb{R}^n$ is a convex and closed set with non-empty interior, and the scalar-valued function $f(\theta, q) : \mathbb{R}^n \times \mathcal{Q} \rightarrow \mathbb{R}$ is convex in the decision vector θ for any $q \in \mathcal{Q} \subseteq \mathbb{R}^l$. The uncertainty q enters into the constraint function $f(\theta, q)$ without assuming any structure, except for the Borel measurability [24] of $f(\theta, \cdot)$ for any fixed θ . In particular, $f(\theta, \cdot)$ may be affected by parametric (possibly nonlinear) and nonparametric uncertainty.

Note that a linear objective function is not essential and the results of the paper still hold for any convex function by a simple relaxation. Specifically, consider a convex objective

function $f_0(\theta)$ and introduce an auxiliary variable t . Then, the optimization in (1) is equivalent to

$$\min_{\theta \in \Theta, t \in \mathbb{R}} t \quad \text{subject to } f_0(\theta) - t \leq 0 \text{ and } f(\theta, q) \leq 0, \forall q \in \mathcal{Q}.$$

Obviously, the above objective function becomes linear in the augmented decision variable (θ, t) and is of the same form as (1). That is, there is no loss of generality to focus on a linear objective function.

B. Motivating Examples

The robust convex optimization in (1) is crucial in many areas of research, see e.g. [5], [12] and references therein for more comprehensive examples. Here we present some important applications for illustration.

Example 1 (Robust MPC). Consider uncertain linear systems

$$x^{k+1} = A(q)x^k + B(q)u^k \quad (2)$$

where $q \in \mathcal{Q}$ represents the system uncertainty. The robust model predictive control (MPC) aims to solve the following optimization problem

$$\begin{aligned} \min_{u^k, \dots, u^{k+h-1}} \max_{q \in \mathcal{Q}} \sum_{j=k}^{k+h-1} g(x^j, u^j) + v(x^{k+h}) \\ \text{subject to } u^j, \dots, u^{k+h-1} \in \mathcal{U} \text{ and (2),} \end{aligned}$$

where g and v are convex functions, and \mathcal{U} is convex and closed. Let $\theta = (u^k, \dots, u^{k+h-1})$, it follows from (2) that the objective function can be rewritten as $J(\theta, q) := \sum_{j=k}^{k+h-1} g(x^j, u^j) + v(x^{k+h})$. Hence, the robust MPC is reformulated as the following RCO

$$\min_{\eta, \theta \in \mathcal{U}^h} \eta \quad \text{subject to } J(\theta, q) - \eta \leq 0, \forall q \in \mathcal{Q}.$$

Example 2 (Distributed robust optimization). Consider the distributed robust optimization problem

$$\min_{\theta \in \Theta} \sum_{j=1}^m f_j(\theta, q_j), \quad (3)$$

where f_j is only known to node j and $q_j \in \mathcal{Q}_j$ represents the uncertainty in node j and its bounding set. Moreover, $f_j(\theta, q_j)$ is convex in θ for any q_j and is Borel measurable in q_j for any fixed θ .

From the worst-case point of view, we are interested in solving the following optimization problem

$$\min_{\theta \in \Theta} \sum_{j=1}^m \left(\max_{q_j \in \mathcal{Q}_j} f_j(\theta, q_j) \right). \quad (4)$$

However, the uncertainty q_j generically enters the objective function $f_j(\theta, q_j)$ in (3) without any specific structure, so that the objective function cannot be explicitly found. To solve (4), we note that it is equivalent to the following optimization problem

$$\min_{\theta \in \Theta, t} \sum_{j=1}^m t_j \quad \text{subject to } \max_{q_j \in \mathcal{Q}_j} f_j(\theta, q_j) - t_j \leq 0, \forall j \in \mathcal{V}. \quad (5)$$

Let $f(t, \theta, q) = [f_1(\theta, q_1) - t_1, \dots, f_m(\theta, q_m) - t_m]'$ where $t = [t_1, \dots, t_m]'$ and $q = [q_1, \dots, q_m]'$ and $\mathcal{Q} = \mathcal{Q}_1 \times \dots \times \mathcal{Q}_m$. Then, the optimization in (5) is equivalent to

$$\min_{\theta \in \Theta, t} \sum_{j=1}^m t_j \quad \text{subject to } f(t, \theta, q) \leq 0, \forall q \in \mathcal{Q}. \quad (6)$$

Clearly, (6) is RCO of the form in (1), except that f_j is only known to node j . However, this is not an issue as discussed in Example 5 in Section III-B.

Example 3 (LASSO). Consider the least squares (LS) problem

$$\min_v \|b - Xv\|,$$

where $X \in \mathbb{R}^{l \times n}$ is the regression matrix and b is the measurement vector. It is well-known that the LS solution has poor numerical properties when the regression matrix is ill-conditioned. A common approach for addressing it is to introduce ℓ^1 regularization technique, which results in a LASSO problem

$$\min_v \{ \|b - Xv\| + \sum_{i=1}^n c_i |v_i| \},$$

where $c_i > 0$ quantifies the robustness of the solution with respect to the i -th column of X . By [25], the LASSO is in fact equivalent to a robust LS problem

$$\min_v \max_{q \in \mathcal{Q}} \|b - (X + q)v\| \quad (7)$$

with the following uncertainty set

$$\mathcal{Q} = \{ [q_1, \dots, q_n] \mid \|q_j\| \leq c_j, j = 1, \dots, n \}.$$

From (7), the LASSO is inherently robust to the uncertainty in the regression matrix X , and the weight factor c_i quantifies its robustness performance. Note that the optimization in (7) can also be reformulated as RCO in (1).

Example 4 (Distribution-free robust optimization). Consider a distribution-free robust optimization under moment constraints

$$\min_{\theta \in \Theta} \max_{q \in \mathcal{P}} \mathbb{E}[f(\theta, q)] \quad (8)$$

where $f(\theta, q)$ is a utility convex function in the decision variable θ for any given realization of the random vector q , and the expectation $\mathbb{E}[\cdot]$ is taken with respect to q . Moreover, \mathcal{P} is a collection of random vectors with the same support, first- and second-moments

$$\mathcal{P} = \{ q : \text{supp}(q) = \mathcal{Q}, \mathbb{E}[q] = \mu, \mathbb{E}[qq'] = \Sigma \}.$$

In light of [26] and the duality theory [27], the optimization problem (8) is equivalent to RCO

$$\min_{\theta, \alpha, \beta, \Omega} \{ \alpha + \mu' \beta + \text{Tr}(\Omega' \Sigma) \}$$

$$\text{subject to } \theta \in \Theta, \alpha + q' \beta + q' \Omega q \geq f(\theta, q), \forall q \in \mathcal{Q}, .$$

Clearly, the optimization (8) is reformulated as RCO of the same form as (1).

Although the stochastic programming (8) is a convex optimization problem, one must often resort to Monte Carlo

sampling to solve it, which is computationally challenging, as it may also need to find an appropriate sampling distribution. Unless f has a special structure, it is very difficult to obtain such a distribution [28]. In the next section, we show how RCO can be effectively solved via a *scenario approach*.

C. Scenario Approach for RCO

The design constraint $f(\theta, q) \leq 0$ for all possible $q \in \mathcal{Q}$ is crucial in the study of robustness of complex systems, e.g. \mathcal{H}_∞ performance of a system affected by the parametric uncertainty and the design of uncertain model predictive control [29]. However, obtaining worst-case solutions has been proved to be computationally difficult, even NP-hard as the uncertainty q may enter into $f(\theta, q)$ in a nonlinear manner. In fact, it is generally very difficult to explicitly characterize the constraint set with uncertainty, i.e.,

$$\{\theta | f(\theta, q) \leq 0, \forall q \in \mathcal{Q}\}, \quad (9)$$

which renders it impossible to directly solve RCO. There are only few cases when the uncertainty set is tractable [12]. Furthermore, this approach introduces undesirable *conservatism*. For these reasons, we adopt the scenario approach.

Instead of satisfying the hard constraint in (9), the idea of this approach is to derive a probabilistic *approximation* by means of a finite number of random constraints, i.e.,

$$\bigcap_{i=1}^{N_{bin}} \{\theta | f(\theta, q^{(i)}) \leq 0\} \quad (10)$$

where N_{bin} is a positive integer representing the constraint size, and $\{q^{(i)}\} \subseteq \mathcal{Q}$ are independent identically distributed (i.i.d.) samples extracted according to an arbitrary absolutely continuous (with respect to the Lebesgue measure) distribution $\mathbb{P}_q(\cdot)$ over \mathcal{Q} .

Regarding the constraint in (10), we only guarantee that most, albeit not all, possible uncertainty constraints in RCO are not violated. Due to the randomness of $\{q^{(i)}\}$, the set of constraint in (10) may be very close to its counterpart (9) in the sense of obtaining a small *violation probability*, which is now formally defined.

Definition 1 (Violation probability). *Given a decision vector $\theta \in \mathbb{R}^n$, the violation probability $V(\theta)$ is defined as*

$$V(\theta) := \mathbb{P}_q\{q \in \mathcal{Q} | f(\theta, q) > 0\}.$$

The multi-sample $q^{1:N_{bin}} := \{q^{(1)}, \dots, q^{(N_{bin})}\}$ is called a *scenario* and the resulting optimization problem under the constraint (10) is referred to as a *scenario problem* (SP)

$$\min_{\theta \in \Theta} c' \theta \quad \text{subject to} \quad f(\theta, q^{(i)}) \leq 0, i = 1, \dots, N_{bin}. \quad (11)$$

In the sequel, let Θ^* be the set of optimal solutions to the SP and Θ_0 be the set of feasible solutions, i.e.,

$$\Theta_0 = \{\theta \in \Theta | f(\theta, q^{(i)}) \leq 0, i = 1, \dots, N_{bin}\}. \quad (12)$$

For the SP, we need the following assumption to study its probabilistic relationship with RCO in (1).

Assumption 1 (Non-empty set of optimal solutions and interior point). *The SP in (11) has a non-empty set of optimal solutions, i.e., $\Theta^* \neq \emptyset$. In addition, there exists a vector $\theta_0 \in \Theta$ such that*

$$f(\theta_0, q^{(i)}) < 0, \forall i = 1, \dots, N_{bin}. \quad (13)$$

The interiority condition (often called Slater's constraint qualification) in (13) implies that there is no duality gap between the primal and dual problems of (11) and the dual problem contains at least an optimal solution [14]. We remark that in robust control it is common to study strict inequalities [29], e.g., when dealing with robust asymptotic stability of a system and therefore this is not a serious restriction. In fact, the set of feasible solutions to (1) is a subset of that of the SP in (11). The main result of the scenario approach for RCO is stated below.

Lemma 1 ([30]). *Assume that there exists a unique solution to (11). Let $\epsilon, \delta \in (0, 1)$, and N_{bin} satisfy the following inequality*

$$\sum_{i=0}^{n-1} \binom{N_{bin}}{i} \epsilon^i (1-\epsilon)^{N_{bin}-i} \leq \delta. \quad (14)$$

Then, with probability at least $1 - \delta$, the solution θ_{sc} of the scenario optimization problem (11) satisfies $V(\theta_{sc}) \leq \epsilon$, i.e.,

$$\mathbb{P}_{q^{1:N_{bin}}} \{V(\theta_{sc}) \leq \epsilon\} \geq 1 - \delta.$$

The uniqueness condition can be relaxed in most cases by introducing a tie-breaking rule, see Section 4.1 of [6]. If the sample complexity N_{bin} satisfies (14), a solution θ_{sc} to (11) approximately solves RCO in (1) with certain probabilistic guarantee. A subsequent problem is to compute the sample complexity, which dictates the smallest number of constraints required in the SP to solve (11). This problem has been addressed in [31] obtaining an improved bound

$$N_{bin} \geq \frac{e}{\epsilon(e-1)} (-\ln \delta + n - 1) \quad (15)$$

where e is the Euler's number. Thus, RCO in (1) can be approximately solved via the SP in (11) with a sufficiently large N_{bin} .

The remaining objective of this paper is to effectively solve the SP in (11) when N_{bin} is large.

III. DISTRIBUTED COMPUTATION SCHEME FOR SCENARIO PROBLEMS

In this section, we introduce a distributed computational framework where many processors (nodes) with limited computational capability are interconnected via a graph. Then, we reformulate the SP in (11) as a distributed optimization problem, which assigns some local constraints to each node and adapts the coupled constraint to the graph structure.

A. Distributed Computing Nodes

Although RCO in (1) can be effectively attacked via the scenario approach, clearly N_{bin} may be large to achieve a high confidence level with small violation probability. For example, in a problem with $n = 32$ variables, setting probability levels

$\epsilon = 0.001$ and $\delta = 10^{-6}$, it follows from (15) that the number of constraints in the SP is $N_{bin} \geq 70898$. For such a large sample complexity N_{bin} , the computational cost for solving the SP (11) becomes very high, which may be far from the computational and memory capacity of a single processor.

To overcome this issue, we propose to use m computing units (nodes) which cooperatively solve the SP in (11) in a distributed fashion. Then, the number of design constraints for node j is reduced to n_j . To maintain the desired probabilistic guarantee, it follows from (15) that $\sum_{j=1}^m n_j \geq N_{bin}$.

A simple heuristic approach is to assign the number of constraints in (11) among nodes proportional to their computing and memory power. In practice, each node can declare the total number of constraints that can be handled. If the number of nodes is comparable to the scenario size N_{bin} , the number of constraints for every node j is significantly reduced, e.g. $n_j \ll N_{bin}$, and n_j can be even as small as one.

The problem is then how to distribute the computational task across multiple nodes to cooperatively solve the SP. To this end, we introduce a directed graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ to model interactions between the computing nodes where $\mathcal{V} := \{1, \dots, m\}$ denotes the set of nodes, and the set of links between nodes is represented by \mathcal{E} . A directed edge $(i, j) \in \mathcal{E}$ exists in the graph if node i directly receives information from node j . Then, the in-neighbors and out-neighbors of node j are respectively defined by $\mathcal{N}_j^{in} = \{i | (j, i) \in \mathcal{E}\}$ and $\mathcal{N}_j^{out} = \{i | (i, j) \in \mathcal{E}\}$. Clearly, every node can directly receive information from its in-neighbors and broadcast information to its out-neighbors. A sequence of directed edges $(i_1, i_2), \dots, (i_{k-1}, i_k)$ with $(i_{j-1}, i_j) \in \mathcal{E}$ for all $j \in \{2, \dots, k\}$ is called a directed path from node i_k to node i_1 . A graph \mathcal{G} is said to contain a *spanning tree* if it has a root node that is connected to any other node in the graph via a directed path, and is *strongly connected* if each node is connected to every other node in the graph via a directed path.

We say that $A = \{a_{ij}\} \in \mathbb{R}^{m \times m}$ is a *row-stochastic* weighting matrix adapted to the underlying graph \mathcal{G} , e.g., $a_{ij} > 0$ if $(i, j) \in \mathcal{E}$ and 0, otherwise, and $a_{jj} = 1 - \sum_{i=1, i \neq j}^m a_{ji} \geq 0$ for all $j \in \mathcal{V}$. Moreover, we denote the associated Laplacian matrix of \mathcal{G} by $\mathcal{L} = I_m - A$. If \mathcal{G} is undirected, A is a symmetric matrix and $\mathcal{N}_j^{in} = \mathcal{N}_j^{out}$, which is simply denoted as \mathcal{N}_j .

Overall, the objective of this paper is to solve the following networked optimization problem.

Problem 1 (Distributed scheme). *Assume that \mathcal{G} is strongly connected. Then, each node computes a solution to the SP in (11) under the following setup:*

- Every node j is able to independently generate n_j i.i.d. samples with an absolutely continuous distribution \mathbb{P}_q , and is not allowed to share these samples with other nodes.
- Every node is able to transmit finite dimensional data per packet via a directed/undirected edge.
- The vector c in the objective function, the constraint function $f(\theta, q)$ and the set Θ are accessible to every node.

In contrast with [8], our approach transmits a fixed dimension state vector among nodes. In addition, each node j only deals with a fixed number n_j of constraints. In [8], each node requires to completely solve *local* SPs under an increasing number of constraints. We provide a more detailed comparison between our approach and [8] in Section IV-C.

B. Reformulation of the Scenario Problem

In this work, we propose recursive algorithms with small computation per iteration to distributedly solve the SP. This is particularly suited when several processors cooperate. The main idea is to introduce “local copies” of θ in each node, and to optimize and update these variables by incrementally learning the constraints until a consensus is reached among all the neighboring nodes. The interactions between nodes are made to (indirectly) obtain the constraint set information from other nodes.

Let $q^{(j1)}, \dots, q^{(jn_j)}$ be the samples that are independently generated in node j according to the distribution \mathbb{P}_q . For simplicity, the local constraint functions are collectively rewritten in a vector form

$$f_j(\theta) := \begin{bmatrix} f(\theta, q^{(j1)}) \\ \vdots \\ f(\theta, q^{(jn_j)}) \end{bmatrix} \in \mathbb{R}^{n_j}.$$

Then, the SP in (11) is equivalent to the following constrained minimization problem

$$\min_{\theta \in \Theta} c' \theta \text{ subject to } f_j(\theta) \leq 0, \forall j \in \mathcal{V}, \quad (16)$$

where $f_j(\theta)$ is only known to node j .

Example 5 (Continuation of Example 2). *In (6), the j -th component function of f is only known to node j . Then, node j can independently extract random samples $\{q_j^{(1)}, \dots, q_j^{(n_j)}\}$ from \mathcal{Q}_j and obtain the local inequality*

$$\tilde{f}_j(\theta, t) := \begin{bmatrix} f_j(\theta, q_j^{(1)}) - t_j \\ \vdots \\ f_j(\theta, q_j^{(n_j)}) - t_j \end{bmatrix} \leq 0, \quad (17)$$

which is only known to node j . Thus, the SP associated with the distributed robust optimization in (6) has the same form of (16) and can be solved as well.

Since each node may have very limited computational and memory capability, the algorithm for each node should be easy to implement with a low computational cost. To achieve this goal, we adopt two different approaches in the sequel for undirected and directed graphs, respectively. The first approach (for undirected graphs) exploits the simple structure of a *primal-dual* sub-gradient algorithm [14] which has an explicit recursive form. Moreover, the interpretation of this approach is natural from the viewpoint of optimization theory. It requires a bidirectional information flow between nodes and therefore it is not applicable to directed graphs. To overcome this limitation, the second approach (for directed graphs) revisits the idea of Polyak random algorithm for convex feasibility problem [32]. We remark that in [32] the algorithms are

centralized and do not address distributed computation, which is resolved in this paper by exploiting the network structure.

Next, we show that the SP can be partially separated by adapting it to the network \mathcal{G} .

Lemma 2 (Optimization equivalence). *Assume that \mathcal{G} contains a spanning tree. Then, the optimal solution to the SP in (11) can be found via the following optimization problem*

$$\begin{aligned} \min_{\theta_1, \dots, \theta_m \in \Theta} \quad & \sum_{j=1}^m c' \theta_j \text{ subject to} \\ & \sum_{i=1}^m a_{ji}(\theta_j - \theta_i) = 0, \\ & f_j(\theta_j) \preceq 0, \forall j \in \mathcal{V}. \end{aligned} \quad (18)$$

Proof. By a slight abuse of notation, let θ be the augmented state of θ_j , i.e., $\theta = [\theta'_1, \dots, \theta'_m]'$, and $\mathcal{L} = I_m - A$, which is the associated Laplacian matrix of the graph \mathcal{G} . Then, the constraint in (18) is compactly written as $(\mathcal{L} \otimes I_n) \theta = 0$. This is equivalent to $\theta_1 = \theta_2 = \dots = \theta_m$ as \mathcal{G} contains a spanning tree [33]. Thus, the above optimization problem is reduced to

$$\min_{\{\theta \in \Theta | f_j(\theta) \preceq 0, \forall j \in \mathcal{V}\}} (m \cdot c' \theta)$$

whose set of optimal solutions is equivalent to that of (11). ■

A nice feature of Lemma 2 is that both the objective function and the constraint in (19) of each node are completely decoupled. The only coupled constraint lies in the consensus constraint in (18), which is required to align the state of each node, and can be handled by exploring the graph structure under *local* interactions. Since each node uses it to learn information from every other node, we need the following assumption.

Assumption 2 (Strong connectivity). *The graph \mathcal{G} is strongly connected.*

As the constraint in (19) is only known to node j , this assumption is clearly necessary. Otherwise, there exists a node i that can never be accessed by some other node j . In this case, it is impossible for node j to find a solution to the SP (11) since the information on $f_i(\theta)$ is always missing to node j .

IV. DISTRIBUTED PRIMAL-DUAL SUB-GRADIENT ALGORITHMS FOR UNDIRECTED GRAPHS

Recently, several papers concentrated on the distributed optimization problem of the form in Lemma 2, see e.g. [11], [15], [16], [34]–[36] and references therein. However, they mostly consider a generic local constraint set, i.e., the local constraint (19) is replaced by $\theta_j \in \Theta_j$ for some convex set Θ_j , rather than having an explicit inequality form. Thus, the proposed algorithms require a projection onto the set Θ_j at each iteration to ensure feasibility. This is easy to perform only if Θ_j has a relatively simple structure, e.g., a half-space or a polyhedron. Unfortunately, the computation of the projection onto the set

$$\Theta_j = \{\theta \in \mathbb{R}^n | f_j(\theta) \preceq 0\} \quad (20)$$

is typically difficult and computational demanding. This work does not use projection to handle the inequality constraints. Rather, we exploit the inequality functions by designing distributed primal-dual algorithms for undirected graphs with the aid of an Lagrangian function. Then, we prove that the recursive algorithm in each node asymptotically converges to some common optimal solution of (11).

Since Θ is closed and convex, the optimization problem in Lemma 2 is reformulated with equality constraints

$$\begin{aligned} \min \quad & \sum_{j=1}^m c' \theta_j + h_\rho(\theta) \\ \text{subject to} \quad & (\mathcal{L}_j \otimes I_n) \theta = 0, g_j(\theta_j) = 0, \forall j \in \mathcal{V} \end{aligned} \quad (21)$$

where \mathcal{L}_j is the j -th row of the Laplacian matrix \mathcal{L} , and $g_j(\theta_j)$ is a function only related to the local constraint of node j , i.e.,

$$g_j(\theta_j) = \begin{bmatrix} d(\theta_j, \Theta) \\ f_j(\theta_j)_+ \end{bmatrix}.$$

The distance function $d(\theta, \Theta)$ measures the distance from the point θ to the set Θ and is obviously convex in θ_j . Since Θ is closed and convex, then $d(\theta, \Theta) = 0$ if and only if $\theta \in \Theta$.

With a slight abuse of notation, we use $\theta = [\theta'_1, \dots, \theta'_m]'$ to denote the augmented state of θ_j . The added quadratic penalty function is defined as

$$h_\rho(\theta) = \frac{\rho}{2} \sum_{j=1}^m \|(\mathcal{L}_j \otimes I_n) \theta\|^2 + \|g_j(\theta_j)\|^2$$

and $\rho > 0$ is a given weighting parameter.

A. Distributed Primal-dual Sub-gradient Algorithm

To solve the optimization problem (21), we focus on the following *Lagrangian*

$$L(\theta, \lambda, \gamma) = \sum_{j=1}^m L_j(\theta, \lambda_j, \gamma_j) \quad (22)$$

with the *local* Lagrangian $L_j(\theta, \lambda_j, \gamma_j)$ defined as

$$L_j = c' \theta_j + \lambda'_j (\mathcal{L}_j \otimes I_n) \theta + \gamma'_j g_j(\theta_j) + h_\rho(\theta)$$

where λ_j and γ_j are the Lagrange multipliers corresponding to (18) and (19), respectively. Then, our objective reduces to find a saddle point $(\theta^*, \lambda^*, \gamma^*)$ of the Lagrangian L in (22), i.e., for any $(\theta, \lambda, \gamma)$, it holds that

$$L(\theta^*, \lambda, \gamma) \leq L(\theta^*, \lambda^*, \gamma^*) \leq L(\theta, \lambda^*, \gamma^*). \quad (23)$$

The existence of a saddle point is ensured under Assumptions 1 and 2, as stated below.

Lemma 3 (Saddle point). *Under Assumptions 1 and 2, there exists a saddle point $(\theta^*, \lambda^*, \gamma^*)$ of the Lagrangian L in (22).*

Proof. Under Assumption 1, it follows from Propositions 5.1.6 and 5.3.1 in [14] that there exists a saddle point for the optimization (11). By the equivalence of the SP in (11) and the problem in Lemma 2, the rest of proof follows. ■

By the Saddle Point Theorem (see e.g. Proposition 5.1.6 in [14]), it is sufficient to find a saddle point of the form (23). In

the section, we design a distributed primal-dual sub-gradient method to achieve this goal.

If $0 \preceq \gamma$, then $L(\theta, \lambda, \gamma)$ is convex in each argument, e.g. $L(\cdot, \lambda, \gamma)$ is convex for any fixed (λ, γ) satisfying $0 \preceq \gamma$. Thus, the following set-valued mappings

$$\begin{aligned} T_j(\theta, \lambda, \gamma) &= \partial_{\theta_j} L(\theta, \lambda, \gamma), \\ P_j(\theta, \lambda, \gamma) &= -\partial_{(\lambda_j, \gamma_j)} L(\theta, \lambda, \gamma) \end{aligned}$$

are well-defined where $\partial_{\theta_j} L(\theta, \lambda, \gamma)$ is the subdifferential of L in θ_j [14]. The optimality of a saddle point $(\theta^*, \lambda^*, \gamma^*)$ becomes $0 \in T_j(\theta^*, \lambda^*, \gamma^*)$ and $0 \in P_j(\theta^*, \lambda^*, \gamma^*)$, which is solved via the following iteration

$$\theta_j^{k+1} = \theta_j^k - \zeta^k \cdot T_j^k \text{ and } \nu_j^{k+1} = \nu_j^k - \zeta^k \cdot P_j^k. \quad (24)$$

Here it is sufficient to arbitrarily select $T_j^k \in T_j(\theta^k, \lambda^k, \gamma^k)$ and $P_j^k \in P_j(\theta^k, \lambda^k, \gamma^k)$. The purpose of ν_j^k is to compute the Lagrange multipliers of $(\lambda_j^*, \gamma_j^*)$. The stepsizes satisfy the following condition

$$\zeta^k > 0, \quad \sum_{k=0}^{\infty} \zeta^k = \infty, \quad \text{and} \quad \sum_{k=0}^{\infty} (\zeta^k)^2 < \infty. \quad (25)$$

Next, we show that the sub-gradient iteration in (24) can be distributedly computed via Algorithm 1 for undirected graphs. For notational simplicity, the dependence of the superscript k , which denotes the number of iterations, is removed. In Algorithm 1, every node keeps updating a triple of state vector and Lagrange multipliers $(\theta_j, \lambda_j, \gamma_j)$ by receiving information only from its neighboring nodes $i \in \mathcal{N}_j$, see Fig. 1. Notice from (22) that (λ_j, γ_j) is a pair of Lagrange multipliers that only appears in the local Lagrangian L_j . This implies that

$$P_j^k = - \begin{bmatrix} \sum_{i=1}^m a_{ji}(\theta_j^k - \theta_i^k) \\ g_j(\theta_j^k) \end{bmatrix}.$$

Clearly, $\sum_{i=1}^m a_{ji}(\theta_j^k - \theta_i^k)$ in P_j^k is computable in node j by receiving information only from in-neighbors of node j . As $g_j(\theta_j^k)$ is a function of local variables, P_j^k is accessible to node j via only local interactions with its in-neighbors. By the additive property of the subdifferential [14], we further obtain from (22) that

$$\begin{aligned} T_j(\theta^k, \lambda^k, \gamma^k) &= c + \sum_{i=1}^m l_{ij} (\lambda_i^k + \rho \cdot (\mathcal{L}_i \otimes I_n) \theta^k) \\ &\quad + s'_j (\gamma_j^k + \rho \cdot g_j(\theta_j^k)), \end{aligned}$$

where l_{ij} is the (i, j) -th element of the Laplacian matrix \mathcal{L} and s_j represents a subgradient of $g_j(\cdot)$ at θ_j , i.e., let ∇_j be a subgradient of $f(\cdot)_+$ at θ_j , then

$$s'_j = \left[\frac{\theta_j - \Pi_{\Theta}(\theta_j)}{\|\theta_j - \Pi_{\Theta}(\theta_j)\|}, \nabla'_j \right] \in \mathbb{R}^{n \times (n_j+1)}. \quad (26)$$

Similarly, the second term in the sum

$$(\mathcal{L}_i \otimes I_n) \theta^k = \sum_{j=1}^m a_{ij} (\theta_i^k - \theta_j^k)$$

is locally computable in node i . Together with the fact \mathcal{G} is undirected, both in-neighbors and out-neighbors of node

Algorithm 1: Distributed primal-dual algorithm for the SP with undirected graphs

- 1: **Initialization:** Each node $j \in \mathcal{V}$ sets $\theta_j = 0$, $\gamma_j = 0$, and $\lambda_j = 0$.
- 2: **Repeat**
- 3: **Local information exchange:** Every node $i \in \mathcal{V}$ broadcasts θ_i to its neighbor $j \in \mathcal{N}_i$, computes $b_i = \sum_{j \in \mathcal{N}_i} a_{ij}(\theta_i - \theta_j)$ after receiving θ_j from neighbor $j \in \mathcal{N}_i$, then broadcasts $\tilde{\lambda}_i = \lambda_i + \rho b_i$ to node $j \in \mathcal{N}_i$, see Fig. 1 for an illustration.
- 4: **Local variables update:** Every node $j \in \mathcal{V}$ updates $(\theta_j, \lambda_j, \gamma_j)$ as follows

$$\lambda_j \leftarrow \lambda_j + \zeta \cdot b_j,$$

$$\gamma_j \leftarrow \gamma_j + \zeta \cdot g_j(\theta_j),$$

$$\theta_j \leftarrow \theta_j - \zeta \cdot \left(c + s'_j \tilde{\gamma}_j + \sum_{i \in \mathcal{N}_j} a_{ij} (\tilde{\lambda}_j - \tilde{\lambda}_i) \right)$$

where $\tilde{\gamma}_j = \gamma_j + \rho \cdot g_j(\theta_j)$, and s_j is a subgradient of $g_j(\cdot)$ at θ_j , see (26).

- 5: **Set** $k = k + 1$.
 - 6: **Until** a predefined stopping rule (e.g., a maximum iteration number) is satisfied.
-

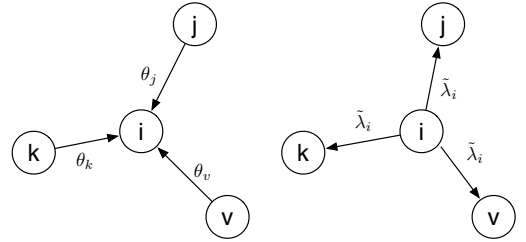


Fig. 1: Local information exchange: every node i receives $\theta_j, \forall j \in \mathcal{N}_i$ from its in-neighbors to compute b_i and $\tilde{\lambda}_i$, after which it broadcasts $\tilde{\lambda}_i$ to out-neighbors.

j are of the same. Thus, the second term in $T_j(\theta^k, \lambda^k, \gamma^k)$ is obtained by aggregating the modified Lagrange multiplier $\tilde{\lambda}_i^k := \lambda_i^k + \rho \cdot (\mathcal{L}_i \otimes I_n) \theta^k$ from its out-neighbors. This further implies that node j is able to compute $T_j(\theta^k, \lambda^k, \gamma^k)$ via local interactions as well.

B. Convergence of Algorithm 1

The update of Lagrange multipliers in Algorithm 1 has interesting interpretation. If θ_j does not satisfy the local constraint, i.e., $\mathbf{1}' g_j(\theta_j) > 0$, some element of the multiplier vector γ_j is strictly increased and a larger penalty is imposed on the augmented Lagrangian L . This forces the update of θ_j to move toward the local feasible set $\Theta \cap \Theta_j$, where Θ_j is given in (20). If γ_j is bounded and the sequence $\{\theta^k\}$ is convergent, it follows that $\sum_{k=1}^{\infty} \zeta^k g_j(\theta_j^k) < \infty$ and $\sup_k \|P_j^k\| < \infty$. In light of (25), this implies that $\liminf_{k \rightarrow \infty} g_j(\theta_j^k) = 0$. Then, the sequence $\{\theta_j^k\}$ will eventually enter the local constraint set $\Theta \cap \Theta_j$. Similarly, the multiplier λ^k will finally drive the state vector θ_j^k to reach a consensus in each node. Based on

these two observations, it follows that $\{\theta_j^k\}$ finally becomes feasible. The convergence of Algorithm 1 is stated and proved.

Theorem 1 (Convergence). *Suppose that Assumptions 1-2 hold and there is a positive r such that $\max\{\|T_j^k\|, \|P_j^k\|\} \leq r$ for all k . Then, the sequence $\{\theta_j^k\}$ of Algorithm 1 with stepsizes given in (25) converges to some common point in the set Θ^* of the optimal solutions to (11).*

Proof. Let $(\theta^*, \lambda^*, \gamma^*)$ be an arbitrary saddle point in Lemma 3. Then, it follows from (24) that

$$\|\theta_j^{k+1} - \theta_j^*\|^2 = \|\theta_j^k - \theta_j^*\|^2 + r^2(\zeta^k)^2 - 2\zeta^k(\theta_j^k - \theta_j^*)'T_j^k$$

Similarly, one can easily obtain

$$\|\nu_j^{k+1} - \nu_j^*\|^2 \leq \|\nu_j^k - \nu_j^*\|^2 + r^2(\zeta^k)^2 - 2\zeta^k(\nu_j^k - \nu_j^*)'P_j^k.$$

For notational simplicity, let

$$z^k = \begin{bmatrix} \theta^k \\ \nu^k \end{bmatrix}, z^* = \begin{bmatrix} \theta^* \\ \nu^* \end{bmatrix}, \text{ and } w^k = \begin{bmatrix} T^k \\ P^k \end{bmatrix}.$$

Then, summing all $j \in \mathcal{V}$ leads to that

$$\|z^{k+1} - z^*\|^2 \leq \|z^k - z^*\|^2 + 2r^2(\zeta^k)^2 - 2\zeta^k(z^k - z^*)'w^k. \quad (27)$$

The rest of the proof is completed by establishing the following two claims.

Claim 1: $(z^k - z^*)'w^k \geq 0$ for all $k \geq 1$.

To show the non-negativeness, we write

$$\begin{aligned} (z^k - z^*)'w^k &= \sum_{j=1}^m \left((c + \sum_{i=1}^m l_{ij}\tilde{\lambda}_i^k + s_j'\tilde{\gamma}_j^k)'(\theta_j^k - \theta_j^*) \right. \\ &\quad \left. - (b_j^k)'(\lambda_j^k - \lambda_j^*) - (g_j(\theta_j^k))'(\gamma_j^k - \gamma_j^*) \right), \end{aligned} \quad (28)$$

where $\tilde{\gamma}_j^k = \gamma_j^k + \rho g_j(\theta_j^k)$ is a modified Lagrange multiplier.

Noting that $g_j(\theta_j^*) = 0$ and $b_i^* = 0$, the sum in (28) is split into four sums. The first sum is the difference between two non-penalized Lagrangians, i.e.,

$$\sum_{j=1}^m (c'\theta_j^k + (\lambda_j^*)'b_j^k + (\gamma_j^*)'g_j(\theta_j^k) - c'\theta_j^*).$$

The second sum involves the Lagrange multiplier λ^k , i.e.,

$$\begin{aligned} &\sum_{j=1}^m \left(\sum_{i=1}^m l_{ij}(\lambda_i^k)'(\theta_j^k - \theta_j^*) - (\lambda_j^k)'b_j^k \right) \\ &= \sum_{i=1}^m (\lambda_i^k)' \sum_{j=1}^m l_{ij}(\theta_j^k - \theta_j^*) - \sum_{j=1}^m (\lambda_j^k)'b_j^k \\ &= \sum_{i=1}^m (\lambda_i^k)'(b_i^k - b_i^*) - \sum_{j=1}^m (\lambda_j^k)'b_j^k \\ &= 0 \end{aligned}$$

where we have used the fact that $b_i^* = 0$ for all $i \in \mathcal{V}$. The third sum involves the Lagrange multiplier γ^k , i.e.,

$$\begin{aligned} &\sum_{j=1}^m (\gamma_j^k)'(s_j^k(\theta_j^k - \theta_j^*) - g_j(\theta_j^k)) \\ &\geq \sum_{j=1}^m (\gamma_j^k)'(g_j(\theta_j^k) - g_j(\theta_j^*) - g_j(\theta_j^k)) \\ &= 0 \end{aligned}$$

where the inequality follows from the fact that $\gamma_j^k \succeq 0$, $g_j(\theta_j^*) = 0$ and s_j^k is a sub-gradient of the vector function $g_j(\theta_j)$ at θ_j^k . The fourth sum involves the penalty term, i.e.,

$$\begin{aligned} &\rho \sum_{j=1}^m \left(\sum_{i=1}^m l_{ij}b_i^k + (s_j^k)'g_j(\theta_j^k) \right)' (\theta_j^k - \theta_j^*) \\ &= \rho \sum_{i=1}^m (b_i^k)'(b_i^k - b_i^*) + g_i(\theta_i^k)'s_i^k(\theta_i^k - \theta_i^*) \\ &\geq \rho \sum_{i=1}^m (\|b_i^k\|^2 + \|g_i(\theta_i^k)\|^2) = 2h_\rho(\theta^k), \end{aligned}$$

where the inequality follows from $b_i^* = 0$, $g_i(\theta_i^*) = 0$ for all $i \in \mathcal{V}$ and the non-negativeness of $g_i(\theta_i^k)$, together with the fact that s_i^k is a sub-gradient of the vector function $g_i(\theta)$ at θ_i^k . Summing the above four sums, we finally obtain that

$$(z^k - z^*)'w^k \geq L(\theta^k, \lambda^*, \gamma^*) - L(\theta^*, \lambda^*, \gamma^*) + h_\rho(\theta^k), \quad (29)$$

which is non-negative by Lemma 3.

Claim 2: $\lim_{k \rightarrow \infty} \theta_j^k = \lim_{k \rightarrow \infty} \theta_i^k \in \Theta^*$ for all $i, j \in \mathcal{V}$.

To this end, jointly with Proposition A.4.4 in [37], (25) and (27), it follows from Claim 1 that the sequence $\{\|z^k - z^*\|\}$ is convergent. Then, $\|z^k\|$ is uniformly bounded. This further implies that the subgradient $\|w^k\|$ is uniformly bounded, e.g., $\|w^k\| \leq \bar{w} < \infty$ for all $k > 0$. By Claim 1 and Proposition A.4.4 in [37], it follows from (27) that

$$\sum_{k=1}^{\infty} \zeta^k (z^k - z^*)'w^k < \infty.$$

Together with (25), we obtain that

$$\liminf_{k \rightarrow \infty} (z^k - z^*)'w^k = 0.$$

In view of (29), it follows that $\liminf_{k \rightarrow \infty} L(\theta^k, \lambda^*, \gamma^*) = L(\theta^*, \lambda^*, \gamma^*)$ and $\liminf_{k \rightarrow \infty} h_\rho(\theta^k) = 0$. Jointly with (22), we finally obtain that $\liminf_{k \rightarrow \infty} \sum_{i=1}^m c'\theta_i^k = \liminf_{k \rightarrow \infty} (\mathbf{1} \otimes c)'\theta^*$ and $\liminf_{k \rightarrow \infty} \theta_i^k = \liminf_{k \rightarrow \infty} \theta_j^k$ for all $i, j \in \mathcal{V}$. That is, there exists an optimal point $\theta_0^* \in \Theta^*$ such that $\liminf_{k \rightarrow \infty} \theta_i^k = \theta_0^*$ for all $i \in \mathcal{V}$. Moreover, one can easily verify that $(\mathbf{1} \otimes \theta_0^*, \lambda^*, \gamma^*)$ is also a saddle point of Lemma 3. Together with Claim 1, it holds that $\{\|\theta_i^k - \theta_0^*\|\}$ converges. Hence, $\lim_{k \rightarrow \infty} \theta_i^k = \theta_0^* \in \Theta^*$ for all $i \in \mathcal{V}$. ■

Corollary 1 (Error bounds). *Under the conditions of Theorem 1, let $\theta^k = \sum_{t=1}^k \zeta^t \theta^t / t^k$ where $t^k = \sum_{t=1}^k \zeta^t$. Then,*

$$\begin{aligned} &L(\check{\theta}^k, \lambda^*, \gamma^*) - L(\theta^*, \lambda^*, \gamma^*) + h_\rho(\check{\theta}^k) \\ &\leq (\|z^1 - z^*\|^2 + 2r^2 \sum_{t=1}^k (\zeta^t)^2) / (2t^k). \end{aligned} \quad (30)$$

Proof. It is straightforward by combining (27) and (29). ■

C. Comparisons with the State-of-the-art

To solve the SP in (11), a distributed setup is proposed in [8] by exchanging the active constraints with neighbors. Specifically, each node j solves a *local* SP of the form

$$\begin{aligned} &\min_{\theta \in \Theta} c'\theta \quad \text{subject to} \\ &f(\theta, q^{(i)}) \leq 0, i \in S_j^k \subseteq \{1, \dots, N_{bin}\} \end{aligned} \quad (31)$$

at each iteration where S_j^0 is the set of indices associated with the random samples generated in node j , and obtains local active constraints, indexed as $ActS_j^k := \{i \in S_j^k | f((\theta_j^k)^*, q^{(i)}) = 0\}$. Here $(\theta_j^k)^*$ is an optimal solution to the local SP in (31), after which it broadcasts its active constraints indexed by $ActS_j^k$ to its out-neighbors. Subsequently, node j updates its local constraint indices as

$$S_j^{k+1} = ActS_j^k \cup (\cup_{i \in \mathcal{N}_j^{in}} ActS_i^k) \cup S_j^0 \quad (32)$$

and returns a local SP of the form (31) replacing S_j^k by S_j^{k+1} . In comparison, one can easily identify several key differences from Algorithm 1.

- (a) Using (31), we cannot guarantee to reduce the computation cost in each node. In particular, it follows from (32) that the number of constraints in each local SP in (31) increases with respect to the number of iterations, and eventually is greater than the total number of active constraints in the SP in (11). In an extreme case, the number of active constraints of (11) can be up to N_{bin} . From this point of view, the computation per iteration in each node is still very demanding. It should be noted that selecting the active constraints of an optimization problem is almost as difficult as solving the entire optimization problem.

In Algorithm 1, it is clear that the computation only requires a few additions and multiplications of vectors, in addition to finding a sub-gradient of a parameterized function $f(\theta, q)$ in θ . It should be noted that the computation of the sub-gradient of $f(\theta, q)$ is unavoidable in almost any optimization algorithm. Clearly, the dimension of γ_j is $n_j + 1$ and $n_j \approx N_{bin}/m$. This implies that the computation cost in each node is greatly reduced as the number of nodes m increases.

- (b) Deciding the active constraints in (31) is very sensitive to the optimal solution $(\theta_j^k)^*$. If $(\theta_j^k)^*$ is not an exact optimal solution or the evaluation of $f((\theta_j^k)^*, q^{(i)})$ is not exact, we cannot correctly identify the index set $ActS_j^k$ of active constraints. In Algorithm 1, there is no such a problem and the local update has certain robustness properties with respect to the round-off errors in computing b_j^k and $g_j(\theta_j^k)$.
- (c) The size of data exchange between nodes in (31) may grow monotonically. Although a quantized index version of (31) is proposed for the channel with bounded communication bandwidth, it needs to compute the vertices of a convex hull per iteration. More importantly, the dimension of the exchanged data per iteration is still larger than that in Algorithm 1.
- (d) The local SP of the form (31) in each node contains several overlapping constraints. Specifically, each constraint set $\{\theta | f(\theta, q^{(i)}) \leq 0\}$ could be handled more than once by every node. This certainly induces redundancy in computation. In Algorithm 1, each inequality is handled exclusively in only one node. From this perspective, Algorithm 1 is of great importance for a node with very limited computational and memory capability.

- (e) It is impossible to describe how the error bounds are reduced with respect to the number of iterations for the distributed algorithms [8].

The primal-dual sub-gradient methods for distributed constrained optimization have been previously used, see e.g., [38]. However, the proposed algorithm originated from the normal Lagrangian (i.e., $\rho = 0$ in (22)). As discussed in [38] after Theorem 1, this usually requires the strict convexity of the Lagrangian to ensure convergence of the primal-dual sequence, which clearly is not satisfied in our case. To remove this strong convexity condition, the authors propose a specially perturbed sub-gradient and assume boundedness on Θ and $\partial_\theta f(\theta, q^{(i)})$. This increases the complexity of the distributed algorithm. In particular, it requires to run up to three consensus algorithms and projects the dual variable onto a bounded ball whose radius must be initially decided, and it is a global parameter. Obviously, Algorithm 1 has a much simpler structure by adopting an augmented Lagrangian in (22), which, to some extent, can be interpreted as the strict convexification of the Lagrangian function. Moreover, the convergence proof of Algorithm 1, which is given in the next subsection, is simpler and easier to understand.

Compared with the distributed alternating direction method of multipliers (ADMM) [39]–[41], the computation of Algorithm 1 is simpler. For example, the ADMM essentially updates the primal sequence as follows

$$\theta^{k+1} \in \operatorname{argmin}_{\theta \in \Theta^m} L_c(\theta, \lambda^k, \gamma^k) \quad (33)$$

where $L_c(\theta, \lambda^k, \gamma^k)$ has a similar form to the augmented Lagrangian $L(\theta, \lambda^k, \gamma^k)$ in (22). That is, it requires to solve an optimization (33) per iteration. In Algorithm 1, we only need to compute one inner iteration to update θ^k by moving along the sub-gradient direction.

D. Extensions to Stochastically time-varying graphs

Algorithms 1 can be easily generalized to the case of stochastically time-varying graphs with a fixed number of nodes. In particular, let the interaction graph at time k be $\mathcal{G}^k := \{\mathcal{V}, \mathcal{E}^k\}$. If $\{\mathcal{G}^k\}$ is an i.i.d. process where the mean graph $\mathbb{E}[\mathcal{G}^k]$ is strongly connected, Theorem 1 continues to hold by following similar lines of proof. For instance, it is easy to show that the SP in (11) is equivalent to

$$\begin{aligned} \min_{\theta_1, \dots, \theta_m \in \Theta} \quad & \sum_{j=1}^m c' \theta_j \quad \text{subject to} \\ & \sum_{i=1}^m \mathbb{E}[a_{ji}^k] (\theta_j - \theta_i) = 0, f(\theta_j, q^{(j)}) \leq 0, \forall j \in \mathcal{V}. \end{aligned} \quad (34)$$

Next, consider a stochastically time-varying augmented Lagrangian

$$L^k(\theta, \lambda, \gamma) = \sum_{j=1}^m L_j^k(\theta, \lambda_j, \gamma_j), \quad (35)$$

where L_j^k is obtained by replacing \mathcal{L}_j with \mathcal{L}_j^k in (22). Moreover, all the elements a_{ij} in Algorithm 1 are replaced by a_{ij}^k . Using the theory of stochastic approximation [42], we

can find a saddle point of $\mathbb{E}[L^k]$, i.e., for any $(\theta, \gamma, \lambda)$, the inequalities

$$\mathbb{E}[L(\theta^*, \lambda, \gamma)] \leq \mathbb{E}[L(\theta^*, \lambda^*, \gamma^*)] \leq \mathbb{E}[L(\theta, \lambda^*, \gamma^*)]$$

hold almost surely. Following a similar reasoning, we can establish the following result, the proof of which is omitted due to the page limitation.

Theorem 2 (Almost sure convergence). *Let Assumption 1 hold and let $\{\mathcal{G}^k\}$ be an i.i.d. sequence with $\mathbb{E}[\mathcal{G}^k]$ strongly connected. If there exists a positive r such that $\max\{\|T_j^k\|, \|P_j^k\|\} \leq r$, the sequence $\{\theta_j^k\}$ of Algorithm 1 with stepsizes in (25) and a_{ij} replaced by a_{ij}^k converges almost surely to some common random point in the set Θ^* of the optimal solutions to (11).*

V. DISTRIBUTED RANDOM PROJECTED ALGORITHMS FOR DIRECTED GRAPHS

In this section, we are concerned with the design of a distributed algorithm for directed graphs. Different from undirected graphs, the information flow between nodes is unidirectional, which results in information unbalance of the network, and renders the primal-dual algorithm inapplicable. To overcome it, we design a consensus algorithm to gather information from in-neighbors and obtain an intermediate state vector. The feasibility is then asymptotically ensured by driving the intermediate state vector toward the local constraint set, which is achieved by updating the solution toward the sub-gradient direction of a randomly selected constraint function. This process is realized by designing a novel distributed variation of a Polyak random algorithm [22], see further comments in Remark 1. The main result is then to prove almost sure convergence of an optimal solution.

A. Distributed Random Projected Algorithm

In Fig. 1, it is clear that the information exchange is bidirectional. In particular, Algorithm 1 requires each node j to use the modified Lagrangian multipliers $\tilde{\lambda}_i$ from its out-neighbors to update the decision vector θ_j . Obviously, this is not implementable for directed graphs, and in this case there is no clear way to design a distributed primal-dual algorithm. For this purpose, we propose a two-stage distributed random projected algorithm

$$v_j^k = \sum_{i=1}^m a_{ji} \theta_i^k - \zeta^k \cdot c, \quad (36)$$

$$\theta_j^{k+1} = \Pi_{\Theta}(v_j^k - \beta \cdot \frac{f(v_j^k, q^{(jw_j^k)})_+}{\|d_j^k\|^2} d_j^k), \quad (37)$$

where $\zeta^k > 0$ is the (deterministic) stepsize given in (25) $\beta \in (0, 2)$ is a constant parameter, $w_j^k \in \{1, \dots, n_j\}$ is a random variable and the vector $d_j^k \in \partial f(v_j^k, q^{(jw_j^k)})_+$ if $f(v_j^k, q^{(jw_j^k)})_+ > 0$ and $d_j^k = d_j$ for some $d_j \neq 0$ if $f(v_j^k, q^{(jw_j^k)})_+ = 0$.

We intuitively explain the key ideas of the above algorithm. The objective of (36) is to distributedly solve an unconstrained

optimization, i.e., the optimization by removing the constraints in (19), see [17] for details. Note that in [17] the double stochasticity of A is required, which is in fact not necessary in our paper. The purpose of (37) is to drive the intermediate state v_j^k toward a randomly selected local constraint set $\Theta \cap \Theta_j^{w_j^k}$, where $\Theta_j^{w_j^k} := \{\theta | f(\theta, q^{(jw_j^k)}) \leq 0\}$. If β is sufficiently small, it is easy to verify (see e.g. [14, Proposition 6.3.1]) that

$$d(\theta_j^{k+1}, \Theta \cap \Theta_j^{w_j^k}) \leq d(v_j^k, \Theta \cap \Theta_j^{w_j^k}).$$

That is, θ_j^{k+1} is closer to the local constraint set $\Theta \cap \Theta_j^{w_j^k}$ than v_j^k . If w_j^k is uniformly selected at random from $\{1, \dots, n_j\}$, we conclude that θ_j^{k+1} is closer to the local constraint set $\Theta \cap \Theta_j$ than v_j^k in the average sense. Once the consensus is achieved among nodes, the state vector θ_j^k in each node asymptotically converges to a point in the feasible set Θ_0 .

Remark 1. *The proposed algorithm is motivated by a generalized Polyak random algorithm [22], which however does not address the distributed design. In this paper, we adapt this algorithm to a directed graph with multiple interconnected nodes and establish its asymptotic optimality for strongly connected digraphs. To the best of our knowledge, the existing work on distributed optimization mostly require the underlying graph to be balanced of the form that the weighting matrix A is doubly stochastic, see e.g. [16]–[19]. Clearly, assuming that the directed graph is balanced is a quite restrictive assumption on the network topology, which is in fact not necessary. This issue has been recently resolved either by combining the gradient descent and the push-sum consensus [20], or augmenting an additional variable for each agent to record the state updates [21]. In comparison, the algorithm in [20] only focuses on the unconstrained optimization, involves nonlinear iterations and requires the updates of four vectors. The algorithm in [21] requires an additional “surplus” vector to record the state update, which increases the computation and communication cost. From this viewpoint, the proposed algorithm of this paper has a simpler structure and is easier to implement, see Algorithm 2 for details.*

B. Convergence of Algorithm 2

To prove convergence, we need the following assumptions, most of which are standard in sub-gradient methods.

Assumption 3 (Randomization and sub-gradient boundedness). *Let the following hold:*

- $\{w_j^k\}$ is an i.i.d. sequence that is uniformly distributed over the set $\{1, \dots, n_j\}$ for any $j \in \mathcal{V}$, and is independent over the index j .
- The sub-gradients d_j^k are uniformly bounded over the set Θ , i.e., there exists a scalar r such that

$$\|d_j^k\| \leq r, \forall j \in \mathcal{V}.$$

Clearly, the designer is free to choose any distribution for drawing the samples w_j^k . Thus, Assumption 3(a) is easy to satisfy. By the property of the sub-gradient and (37), a sufficient condition for Assumption 3(b) is that Θ is bounded.

Algorithm 2: Distributed random projection algorithm for the SP with directed graphs

- 1: **Initialization:** For each node $j \in \mathcal{V}$ set $\theta_j = 0$.
 - 2: **Repeat**
 - 3: **Local information exchange:** Every node $j \in \mathcal{V}$ broadcasts θ_j to its out-neighboring nodes.
 - 4: **Local variables update:** Every node $j \in \mathcal{V}$ receives the state vector θ_i from its in-neighbor $i \in \mathcal{N}_j^{\text{in}}$ and updates it as follows
 - $v_j = \sum_{i \in \mathcal{N}_j^{\text{in}}} a_{ji} \theta_i - \zeta c$ where the stepsize ζ is given in (25).
 - Draw $w_j \in \{1, \dots, n_j\}$ uniformly at random.
 - $\theta_j \leftarrow \Pi_{\Theta}(v_j - \beta \cdot \frac{f(v_j, q^{(jw_j)})_+}{\|d_j\|^2} d_j)$ where d_j is defined in (37).
 - 5: **Set** $k = k + 1$.
 - 6: **Until** a predefined stopping rule (e.g., a maximum iteration number) is satisfied.
-

We now present the convergence result on the distributed random algorithm.

Theorem 3 (Almost sure convergence). *Suppose that Assumptions 1-3 hold. The sequence $\{\theta_j^k\}$ of Algorithm 2 converges almost surely to some common point in the set Θ^* of the optimal solutions to (11).*

C. Proof of Theorem 3

The proof is roughly divided into three parts. The first part establishes a stochastically “decreasing” result, see Lemma 4. That is, the distance of θ^{k+1} to some optimal point θ^* is “stochastically” closer than that of θ^k . The second part essentially shows the asymptotic feasibility of the state vector θ_j^k , see Lemma 5. Finally, the last part establishes an asymptotic consensus result in Lemma 7, which shows that the sequence $\{\theta_j^k\}$ converge to some common value for all $j \in \mathcal{V}$. Combining these results, we show that $\{\theta_j^k\}$ converges almost surely to some common random point in the set Θ^* .

Now, we establish a stochastically “decreasing” result.

Lemma 4 (Stochastically decreasing). *Let \mathcal{F}^k be the sigma-field generated by the random variables $\{w_j^t, j \in \mathcal{V}\}$ up to time k , i.e.,*

$$\mathcal{F}^k = \{w^0, \dots, w^k\} \quad (38)$$

and $\hat{\theta}_j^k = \sum_{i=1}^m a_{ji} \theta_i^k$, where θ_i^k is generated in Algorithm 2.

Under Assumptions 1 and 3, it holds almost surely that for all $j \in \mathcal{V}$ and $k \geq \tilde{k}$, which is a sufficiently large number,

$$\begin{aligned} \mathbb{E}[\|\theta_j^{k+1} - \theta^*\|^2 | \mathcal{F}_k] &\leq (1 + r_1(\zeta^k)^2) \|\hat{\theta}_j^k - \theta^*\|^2 \\ &\quad - 2\zeta^k c'(y_j^k - \theta^*) - r_2(\|\hat{\theta}_j^k - y_j^k\|^2) + r_3(\zeta^k)^2, \end{aligned} \quad (39)$$

where $r_i > 0, i \in \{1, 2, 3\}$, $\theta^* \in \Theta^*$ and $y_j^k = \Pi_{\Theta_0}(\hat{\theta}_j^k)$ with Θ_0 given in (12).

Proof. The proof mostly follows from [22], which however only focuses on the centralized version of Algorithm 2. By the comments after Assumption 2 of [22], it is clear that all

conditions in [22, Proposition 1] are satisfied. By the row stochasticity of A , i.e., $\sum_{i=1}^m a_{ji} = 1$, it follows that (36) can be also written as

$$v_j^k = \hat{\theta}_j^k - \zeta^k \cdot \nabla \left(c' \hat{\theta}_j^k \right),$$

where $\nabla \left(c' \hat{\theta}_j^k \right)$ is a gradient of the linear function $c'\theta$ evaluated at $\theta = \hat{\theta}_j^k$. The rest of proof is trivial by replacing x_{k-1} in (21) of [22] with $\hat{\theta}_j^k$. The details are omitted. ■

The second result essentially ensures the local feasibility.

Lemma 5 (Feasibility guarantee). *Let y_j^k be given in Lemma 4. If $\lim_{k \rightarrow \infty} \|v_j^k - y_j^k\| = 0$, it holds $\lim_{k \rightarrow \infty} \|\theta_j^{k+1} - y_j^k\| = 0$ for any $j \in \mathcal{V}$.*

Proof. Since $f(y_j^k, q^{(jw_j^k)})_+ = 0$, it follows from Lemma 1 of [22] that

$$\|\theta_j^{k+1} - y_j^k\|^2 \leq \|v_j^k - y_j^k\|^2 - \beta(2 - \beta) \frac{(f(v_j^k, q^{(jw_j^k)})_+)^2}{\|d_j^k\|^2}.$$

Together with the fact that $\beta \in (0, 2)$, then $\|\theta_j^{k+1} - y_j^k\|^2 \leq \|v_j^k - y_j^k\|^2$. Taking limits on both sides, the result follows. ■

Finally, we prove an asymptotic consensus result under Assumption 2 where the consensus value is a weighted average of the state vector in each node. This is different than the case of balanced graphs. For a strongly connected digraph \mathcal{G} , we have some preliminary results on its weighting matrix A by directly using the Perron Theorem [43].

Lemma 6 (Left eigenvector). *Under Assumption 2, there exists a normalized left eigenvector $\pi \in \mathbb{R}^m$ of A such that*

$$\pi' A = \pi', \sum_{j=1}^m \pi_j = 1 \text{ and } \pi_j > 0, \forall j \in \mathcal{V}. \quad (40)$$

Moreover, the spectral radius of the row-stochastic matrix $A - \mathbf{1}\pi'$ is strictly less than one.

Lemma 7 (Asymptotic consensus). *Consider the following iteration*

$$\theta_j^{k+1} = \sum_{i=1}^m a_{ji} \theta_i^k + n_j^k, \forall j \in \mathcal{V}.$$

Suppose that \mathcal{G} is strongly connected and $\lim_{k \rightarrow \infty} \|n_j^k\| = 0$. Let $\bar{\theta}^k = \sum_{i=1}^m \pi_i \theta_i^k$, where π_i is given in (40), it holds that

$$\lim_{k \rightarrow \infty} \|\theta_j^k - \bar{\theta}^k\| = 0, \forall j \in \mathcal{V}. \quad (41)$$

Proof. Clearly, we can compactly write $\bar{\theta}^k = (\pi' \otimes I_n) \theta^k$. In view of (40) and (41), we have the following relation

$$\mathbf{1}(\pi' \otimes I_n) \theta^{k+1} = \mathbf{1}(\pi' \otimes I_n) \theta^k + \mathbf{1}(\pi' \otimes I_n) n^k. \quad (42)$$

Let $\delta^k = ((I_n - \mathbf{1}\pi') \otimes I_n) \theta^k$, which is a vector of displacement from the weighted average. Then, it follows from (42) that

$$\delta^{k+1} = ((A - \mathbf{1}\pi') \otimes I_n) \delta^k + ((I - \mathbf{1}\pi') \otimes I_n) n^k.$$

Define ρ as the spectral radius of $(A - \mathbf{1}\pi') \otimes I_n$, it is clear from Lemma 6 that $0 < \rho < 1$. Jointly with the fact

that $\lim_{k \rightarrow \infty} \|n^k\| = 0$ and Lemma 6.1.1 [24], it follows that $\lim_{k \rightarrow \infty} \|\delta^k\| = 0$. ■

The proof also depends crucially on the well-known supermartingale convergence theorem, which is due to Robbins-Siegmund [44], see also Proposition A.4.5 in [37]. This result is now restated for completeness.

Theorem 4 (Super-martingale convergence theorem). *Let $\{y^k\}, \{z^k\}, \{w^k\}$ and $\{v^k\}$ be four non-negative sequences of random variables, and let $\mathcal{F}^k, k = 0, 1, \dots$, be sets of random variables such that $\mathcal{F}^k \subseteq \mathcal{F}^{k+1}$ for all k . Assume that*

- (a) *For each k , let y^k, z^k, w^k and v^k be functions of the random variables in \mathcal{F}^k .*
- (b) *The inequalities hold almost surely*

$$\mathbb{E}[y^{k+1} | \mathcal{F}^k] \leq (1 + v^k)y^k - z^k + w^k, k = 0, 1, \dots, \text{ and}$$

$$\sum_{k=0}^{\infty} w^k < \infty, \quad \sum_{k=0}^{\infty} v^k < \infty.$$

Then, $\{y^k\}$ converges almost surely to a nonnegative random variable y , and $\sum_{k=0}^{\infty} z^k < \infty$.

Combine the above, we are ready to prove Theorem 3.

Proof of Theorem 3. By the convexity of $\|\cdot\|^2$ and the row stochasticity of A , i.e., $\sum_{i=1}^m a_{ji} = 1$, it follows that

$$\|\hat{\theta}_j^k - \theta^*\|^2 \leq \sum_{i=1}^m a_{ji} \|\theta_i^k - \theta^*\|^2.$$

Jointly with (39), we obtain that for all $k \geq \tilde{k}$,

$$\begin{aligned} \mathbb{E}[\|\theta_j^{k+1} - \theta^*\|^2 | \mathcal{F}_k] &\leq (1 + r_1(\zeta^k)^2) \sum_{i=1}^m a_{ji} \|\theta_i^k - \theta^*\|^2 \\ &\quad - 2\zeta^k c'(y_j^k - \theta^*) - r_2(\|\hat{\theta}_j^k - y_j^k\|^2) + r_3(\zeta^k)^2, \end{aligned} \quad (43)$$

where the sigma-field \mathcal{F}^k is given in (38).

Under Assumption 2, the weighting matrix A of \mathcal{G} is only row stochastic, and not doubly stochastic, which is assumed in [16]. This implies that the first term in (36) does not satisfy average consensus. Instead, it converges to the weighted average consensus where the weight is determined by the left eigenvector $\pi \in \mathbb{R}^m$ of A associated with the simple eigenvalue 1, i.e., $\pi' A = \pi$, see Lemma 7. Since the graph \mathcal{G} is strongly connected, it is clear that $\pi_j > 0$ for all $j \in \mathcal{V}$.

Then, we multiply both sides of (43) with π_j and sum over j , which leads to

$$\begin{aligned} &\mathbb{E}\left[\sum_{j=1}^m \pi_j \|\theta_j^{k+1} - \theta^*\|^2 | \mathcal{F}_k\right] \\ &\leq (1 + r_1(\zeta^k)^2) \sum_{j=1}^m \pi_j \left(\sum_{i=1}^m a_{ji} \|\theta_i^k - \theta^*\|^2\right) \\ &\quad - 2\zeta^k c'(\bar{y}^k - \theta^*) - \sum_{j=1}^m \pi_j \left(r_2(\|\hat{\theta}_j^k - y_j^k\|^2) + r_3(\zeta^k)^2\right) \\ &\leq (1 + r_1(\zeta^k)^2) \sum_{j=1}^m \pi_j \|\theta_j^k - \theta^*\|^2 \\ &\quad - 2\zeta^k c'(\bar{y}^k - \theta^*) - r_2 \sum_{j=1}^m \pi_j (\|\hat{\theta}_j^k - y_j^k\|^2) + r_3(\zeta^k)^2 \end{aligned} \quad (44)$$

where the first inequality uses the fact that $\bar{y}^k = \sum_{j=1}^m \pi_j y_j^k$ and $\sum_{j=1}^m \pi_j = 1$. The second inequality follows from the definition of π , i.e., $\pi_j = \sum_{i=1}^m \pi_i a_{ji}$.

By Theorem 4, it holds almost surely that $\{\sum_{j=1}^m \pi_j \|\theta_j^k - \theta^*\|^2\}$ is convergent for any $j \in \mathcal{V}$ and $\theta^* \in \Theta^*$,

$$\sum_{k=1}^{\infty} \zeta^k c'(\bar{y}^k - \theta^*) < \infty \quad (45)$$

and

$$\sum_{k=1}^{\infty} \sum_{j=1}^m \pi_j \|\hat{\theta}_j^k - y_j^k\|^2 < \infty. \quad (46)$$

The rest of the proof is completed by showing the following two claims.

Claim 1: $\{\|\bar{y}^k - \theta^*\|\}$ converges almost surely.

In light of (46), it holds that $\{\|y_j^k - \hat{\theta}_j^k\|\}$ converges to zero almost surely. Since $\zeta^k \rightarrow 0$, it follows from (36) that $\{\|v_j^k - \hat{\theta}_j^k\|\}$ converges almost surely to zero as well. Combing the preceding two relations, it holds almost surely that $\lim_{k \rightarrow \infty} \|y_j^k - v_j^k\| = 0$. Together with Lemma 5, it holds almost surely that $\lim_{k \rightarrow \infty} \|\theta_j^{k+1} - y_j^k\| = 0$ for any $j \in \mathcal{V}$. Since $\{\sum_{j=1}^m \pi_j \|\theta_j^{k+1} - \theta^*\|^2\}$ converges almost surely, this implies that $\{\sum_{j=1}^m \pi_j \|y_j^k - \theta^*\|^2\}$ converges as well.

By (36) and (37), we have the following dynamics

$$\theta_j^{k+1} = \sum_{i=1}^m a_{ji} \theta_i^k + n_j^k \quad (47)$$

where $n_j^k = \theta_j^{k+1} - v_j^k - \zeta^k c$. Since the inequality

$$\|n_j^k\| \leq \|\theta_j^{k+1} - y_j^k\| + \|y_j^k - v_j^k\| + \zeta^k \|c\|,$$

holds, it is obvious that $\lim_{k \rightarrow \infty} \|n_j^k\| = 0$ almost surely. Together with Lemma 7, we obtain that $\lim_{k \rightarrow \infty} \|\theta_j^k - \bar{\theta}^k\| = 0$ almost surely.

Since $\pi_j = \sum_{i=1}^m a_{ji} \pi_i$, it holds that $\bar{\theta}^k = \sum_{j=1}^m \pi_j \theta_j^k = \sum_{i=1}^m \pi_i \hat{\theta}_i^k$. Then, we obtain that

$$\begin{aligned} \|\hat{\theta}_j^k - \sum_{i=1}^m \pi_i \hat{\theta}_i^k\| &= \left\| \sum_{i=1}^m a_{ji} \theta_i - \bar{\theta}^k \right\| \\ &\leq \sum_{i=1}^m a_{ji} \|\theta_i^k - \bar{\theta}^k\| \rightarrow 0 \text{ as } k \rightarrow \infty. \end{aligned}$$

Since $\lim_{k \rightarrow \infty} \|y_j^k - \hat{\theta}_j^k\| = 0$, it follows that $\|y_j^k - \bar{y}^k\| \leq \|y_j^k - \hat{\theta}_j^k\| + \|\hat{\theta}_j^k - \sum_{i=1}^m \pi_i \hat{\theta}_i^k\| + \sum_{i=1}^m \pi_i \|\hat{\theta}_i^k - y_i^k\|$, which converges almost surely to zero as $k \rightarrow \infty$ by using the above relations. Jointly with the fact that $\{\sum_{j=1}^m \pi_j \|y_j^k - \theta^*\|^2\}$ converges, we obtain that $\{\|\bar{y}^k - \theta^*\|\}$ converges almost surely.

Claim 2: There exists $\theta_0^* \in \Theta^*$ such that $\lim_{k \rightarrow \infty} \theta_j^k = \theta_0^*$ for all $j \in \mathcal{V}$ with probability one.

By (25) and (45), it follows that $\liminf_{k \rightarrow \infty} c' \bar{y}^k = c' \theta^*$, which implies that there exists a subsequence of $\{\bar{y}^k\}$ that converges almost surely to some point in the optimal set Θ^* , which is denoted as θ_0^* . Jointly with Claim 1 that $\{\|\bar{y}^k - \theta_0^*\|\}$ converges, it follows that $\lim_{k \rightarrow \infty} \bar{y}^k = \theta_0^*$ almost surely. Finally, we note that $\|\theta_j^{k+1} - \theta_0^*\| \leq \|\theta_j^{k+1} - y_j^{k+1}\| + \|y_j^{k+1} - \bar{y}^{k+1}\| + \|\bar{y}^{k+1} - \theta_0^*\|$, which converges almost surely to zero as $k \rightarrow \infty$. Thus, Claim 2 is proved. ■

Corollary 2 (Error bounds). *Under the conditions of Theorem 3, let $\tilde{y}^k = \frac{1}{t^k} \sum_{t=1}^k \zeta^t y^t$ and $e^k = c'(\tilde{y}^k - \theta^*)$. Then, for all $k \geq \tilde{k}$, it holds that*

$$0 \leq \mathbb{E}[e^k] \leq \frac{c^k}{2t^k} \text{ and } \mathbb{E}[\|\theta_j^k - y_j^k\|^2] \leq \frac{c^k}{a_{jj}\pi_j k} \quad (48)$$

where $y_j^k = \Pi_{\Theta_0}(\hat{\theta}_j^k)$ is feasible and

$$c^k = \exp(r_1 \sum_{t=1}^k (\zeta^t)^2) \left(\sum_{j=1}^m \pi_j \|\theta_j^1 - \theta^*\|^2 + r_3 \sum_{t=1}^k (\zeta^t)^2 \right).$$

Proof. Note that y_j^k is feasible and $\prod_{t=1}^k (1 + r_1(\zeta^t)^2) \leq \prod_{t=1}^k \exp(r_1(\zeta^t)^2) < \infty$. By (44), the proof requires tedious but easy algebraic operations and is omitted to save space. ■

As in Section IV-D, Algorithm 2 can also be modified to deal with the case of stochastically time-varying graphs.

D. Comparison with the Distributed Primal-Dual Algorithm

In this subsection, we compare the previously two algorithms. First, although both algorithms are designed from different perspectives, they essentially converge as fast as $O(1/\sum_{t=1}^k \zeta^t)$. Let $0 < \alpha \leq 0.5$, it follows from (25) that it suffices to select $\zeta^t = t^{-(0.5+\alpha)}$, and

$$\sum_{t=1}^k \zeta^t \approx \int_0^k t^{-(0.5+\alpha)} dt = \begin{cases} k^{0.5-\alpha}, & \text{if } 0 < \alpha < 0.5, \\ \ln k, & \text{if } \alpha = 0.5. \end{cases}$$

This implies that the convergence rate of both algorithms can be as fast as $O(1/\sqrt{k})$, which is an optimal rate for a generic sub-gradient algorithm, see Page 9 in [45].

Second, the primal-dual algorithm is originated from sub-gradient methods for finding a saddle point of the augmented Lagrangian. In [37], there are quite a few methods to accelerate the sub-gradient method, which may provide many opportunities to accelerate the networked primal-dual algorithms. This is not obvious for Algorithm 2 since there is no clear way to accelerate its convergence.

Third, the computational cost of both algorithms is low at each iteration. The algorithms are well-suited for the computing nodes with limited computation and memory capability.

VI. APPLICATION EXAMPLE: ROBUST IDENTIFICATION

To illustrate effectiveness of the proposed distributed algorithms, we consider a RCO problem in (1) with linearly structured uncertainties in an identification problem where we seek to estimate the impulse response θ of a discrete-time system for its input u and output y .

Assume that the system is linear, single input single output and of order n , and that u is zero for negative time indices and θ, u and y are related by the convolution equations $y = U\theta$ where U is a lower-triangular Toeplitz matrix whose first column is u , i.e., let $u = [u_1, \dots, u_n]'$, then

$$U = \begin{bmatrix} u_1 & 0 & \dots & 0 \\ u_2 & u_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ u_n & u_{n-1} & u_2 & u_1 \end{bmatrix}.$$

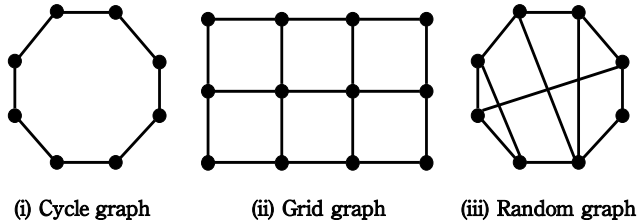


Fig. 2: Three types of graphs.

Suppose that the actual input and output are $u + \delta u$ and $y + \delta y$, respectively. Then, the standard least squares (LS) are not appropriate as the perturbation δu and δy are unknown. To solve it, let $q = [\delta u', \delta y']' \in \mathbb{R}^{2n}$. From the worst point of view, θ is obtained by solving a RCO problem

$$\min_{\theta, t} t \text{ subject to } \|(y + \delta y) - (U + \delta U)\theta\| \leq t, \forall q \in \mathcal{Q}. \quad (49)$$

If $\mathcal{Q} = \{q | \|q\|_\infty \leq \rho\}$, where ρ represents the uncertainty size, it is a structured robust LS problem, which is NP-complete [46]. Thus, we approximately solve it by the scenario approach via distributed Algorithms 1 and 2, and set $u = [1 \ 2 \ 3]'$ and $y = [4 \ 5 \ 6]'$. While for the SP, we consider $\epsilon = 0.002$ and $\delta = 10^{-4}$. This implies from (15) that $N_{bin} \geq 8868$. Here we set $N_{bin} = 10000$ and each node independently extracts samples via a uniform distribution over the uncertainty set \mathcal{Q} .

We adopt three types of undirected graphs, see Fig. 2 where the random graph is obtained by further connecting node i to node $j (\neq i + 1)$ with probability $p = 0.2$ in a cycle graph. A directed random graph is originated from the undirected one. Specifically, node i is connected to node $i + 1$ in the clockwise direction, and the direction of every other link is randomly selected with equal probability.

Given an uncertainty size, define the maximum of the scenario-based residuals by

$$r(\theta, \rho) = \max_{i=1, \dots, N_{bin}} \|(y + \delta y^{(i)}) - (U + \delta U^{(i)})\theta\|. \quad (50)$$

Let $\theta_{ls} = U^{-1}y$ be the solution of the standard LS and θ_{sc} be solution to the SP of (49), which is computed by Algorithm 1. We depict the maximum residuals of (50) in Fig. 3a under different sizes of uncertainty, which shows the robustness of the solution of the SP. Then, we compare the convergence behavior of the proposed algorithms for the SP of (49) with $N_{bin} = 10000, \rho = 0.2$ and $\zeta^k = 2/k$ in (25). Fig. 3b shows that Algorithm 1 converges to a solution of the SP much faster than that of Algorithm 2.

Since for both algorithms the dimensions of the data in crossing a communication link and being stored and retrieved in local memory are constant, one can argue that the total time to run our algorithms is essentially given by

$$T_{total} = T_{comp} + \alpha \cdot N_{iter},$$

where T_{comp} is the time attributed just to computation, α is a constant which mainly depends on the network topology, the communication protocol and the memory access speed, and N_{iter} denotes the number of iterations. Let T_i^k be the time

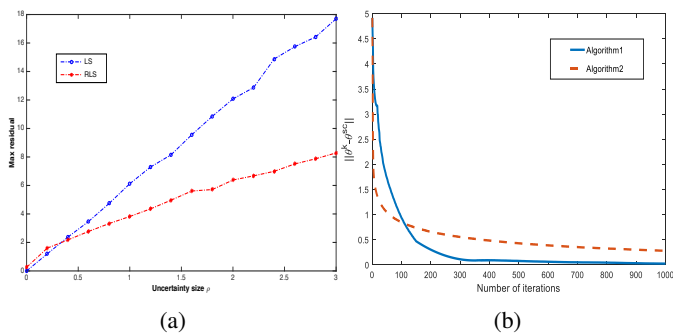


Fig. 3: (a) Maximum residual versus uncertainty size. (b) Convergence behaviors of Algorithms 1 and 2 with $\beta = 1.5$ on undirected and directed random graphs with $m = 100$.

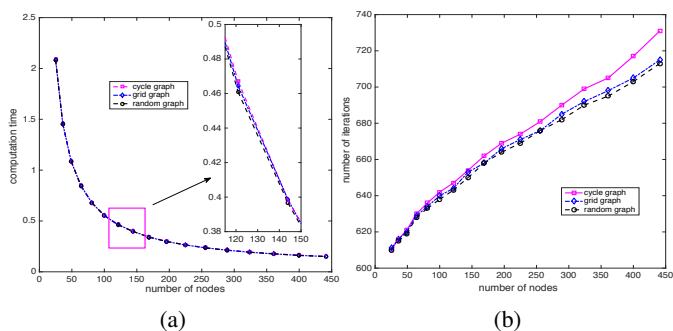


Fig. 4: Performance of Algorithm 1 over three types of network topologies. (a) The time (second) attributed just to computation versus the network size. (b) The number of iterations versus the network size.

cost to compute the k -th iteration of node i , i.e., running steps 3 and 4 in Algorithm 1. Then, it follows from [47, Section 1.2.2] that $T_{comp} = \sum_k \max_{i \in \mathcal{V}} \{T_i^k\}$. Fig. 4a illustrates how the number of nodes affects T_{comp} , which decreases rapidly if the node number m is small, and is indistinguishable for three types of network topologies as each node only involves simple numerical operations. This is consistent with our objective to reduce the computation cost of each node. Moreover, Fig. 4a also indicates that T_{comp}/m is uniformly bounded away from zero, showing the practicability of the proposed distributed algorithm [47, Section 1.2.2]. Ideally, T_{comp}/m needs to be a constant, which is however not attainable [47, Section 1.2.2].

Fig. 4b illustrates that the graph with denser communication links requires a smaller number of iterations, which is clearly consistent with our intuition as the information is mixing faster over a denser graph. However, this requires a higher communication cost. By Fig. 4, one can conclude that designing an optimal topology is extremely complicated, and requires an optimal tradeoff among the communication topology, the number of nodes, and the computation and storage capacity of a single node, some of which are highly coupled. Similar phenomenon can be observed for Algorithm 2 and is not included to save space.

VII. CONCLUSION

In this work, we developed distributed algorithms to collaboratively solve RCO via the SP, which possibly has a large number of constraints. Two distributed algorithms with very simple structure were provided for undirected and directed graphs, respectively. Compared with the existing results, the complexity per iteration of the proposed algorithms is significantly reduced. Future work will focus on exploiting the structure of the parametrized constraint functions to reduce the computation cost.

ACKNOWLEDGEMENT

The authors would like to thank the Associate Editor and anonymous reviewers for their very constructive comments, which greatly improve the quality of this work.

REFERENCES

- [1] G. C. Calafiore and M. C. Campi, "The scenario approach to robust control design," *IEEE Transactions on Automatic Control*, vol. 51, no. 5, pp. 742–753, 2006.
- [2] R. Tempo, G. Calafiore, and F. Dabbene, *Randomized Algorithms for Analysis and Control of Uncertain Systems, with Applications*. Springer-Verlag London, 2013.
- [3] A. Ben-Tal and A. Nemirovski, "Robust convex optimization," *Mathematics of Operations Research*, vol. 23, no. 4, pp. 769–805, 1998.
- [4] C. Scherer, "Relaxations for robust linear matrix inequality problems with verifications for exactness," *SIAM Journal on Matrix Analysis and Applications*, vol. 27, no. 2, pp. 365–395, 2005.
- [5] D. Bertsimas, D. B. Brown, and C. Caramanis, "Theory and applications of robust optimization," *SIAM Review*, vol. 53, no. 3, pp. 464–501, 2011.
- [6] G. C. Calafiore and M. C. Campi, "Uncertain convex programs: randomized solutions and confidence levels," *Mathematical Programming*, vol. 102, pp. 25–46, 2005.
- [7] G. C. Calafiore, F. Dabbene, and R. Tempo, "Research on probabilistic methods for control system design," *Automatica*, vol. 47, no. 7, pp. 1279–1293, 2011.
- [8] L. Carlone, V. Srivastava, F. Bullo, and G. C. Calafiore, "Distributed random convex programming via constraints consensus," *SIAM Journal on Control and Optimization*, vol. 52, no. 1, pp. 629–662, 2014.
- [9] G. Notarstefano and F. Bullo, "Distributed abstract optimization via constraints consensus: Theory and applications," *IEEE Transactions on Automatic Control*, vol. 56, no. 10, pp. 2247–2261, 2011.
- [10] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [11] A. Nedich, "Convergence rate of distributed averaging dynamics and optimization in networks," *Foundations and Trends® in Systems and Control*, vol. 2, no. 1, pp. 1–100, 2015.
- [12] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski, *Robust Optimization*. Princeton University Press, 2009.
- [13] B. L. Gorissen, İ. Yamkoğlu, and D. den Hertog, "A practical guide to robust optimization," *Omega*, vol. 53, pp. 124–137, 2015.
- [14] D. P. Bertsekas, *Nonlinear Programming*. Athena Scientific, 1999.
- [15] N. Chatzipanagiotis, D. Dentecheva, and M. M. Zavlanos, "An augmented Lagrangian method for distributed optimization," *Mathematical Programming*, vol. 152, no. 1-2, pp. 405–434, 2015.
- [16] S. Lee and A. Nedich, "Asynchronous gossip-based random projection algorithms over networks," *IEEE Transactions on Automatic Control*, vol. 61, no. 4, pp. 953–968, 2016.
- [17] A. Nedich and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [18] J. C. Duchi, A. Agarwal, and M. J. Wainwright, "Dual averaging for distributed optimization: convergence analysis and network scaling," *IEEE Transactions on Automatic Control*, vol. 57, no. 3, pp. 592–606, 2012.
- [19] B. Ghahesifard and J. Cortés, "Distributed continuous-time convex optimization on weight-balanced digraphs," *IEEE Transactions on Automatic Control*, vol. 59, no. 3, pp. 781–786, 2014.

- [20] A. Nedich and A. Olshevsky, "Distributed optimization over time-varying directed graphs," *IEEE Transactions on Automatic Control*, vol. 60, no. 3, pp. 601–615, 2015.
- [21] C. Xi and U. A. Khan, "Directed-distributed gradient descent," in *53rd Annual Allerton Conference on Communication, Control, and Computing, Monticello, IL, USA*, 2015.
- [22] A. Nedich, "Random algorithms for convex minimization problems," *Mathematical Programming*, vol. 129, no. 2, pp. 225–253, 2011.
- [23] K. You and R. Tempo, "Parallel computation for robust convex programs over networks," in *American Control Conference, Boston, US*, 2016.
- [24] R. Ash and C. Doléans-Dade, *Probability and Measure Theory*. Academic Press, 2000.
- [25] H. Xu, C. Caramanis, and S. Mannor, "Robust regression and LASSO," *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3561–3574, 2010.
- [26] E. Delage and Y. Ye, "Distributionally robust optimization under moment uncertainty with application to data-driven problems," *Operations Research*, vol. 58, no. 3, pp. 595–612, 2010.
- [27] A. Shapiro, "On duality theory of conic linear problems," in *Semi-infinite Programming*. Springer, 2001, pp. 135–165.
- [28] B. Barmish and C. M. Lagoa, "The uniform distribution: A rigorous justification for its use in robustness analysis," *Mathematics of Control, Signals and Systems*, vol. 10, no. 3, pp. 203–222, 1997.
- [29] I. R. Petersen and R. Tempo, "Robust control of uncertain systems: classical results and recent developments," *Automatica*, vol. 50, pp. 1315–1335, 2014.
- [30] M. C. Campi and S. Garatti, "The exact feasibility of randomized solutions of uncertain convex programs," *SIAM Journal on Optimization*, vol. 19, no. 3, pp. 1211–1230, 2008.
- [31] T. Alamo, R. Tempo, A. Luque, and D. R. Ramirez, "Randomized methods for design of uncertain systems: Sample complexity and sequential algorithms," *Automatica*, vol. 52, pp. 160–172, 2015.
- [32] B. Polyak, "Random algorithms for solving convex inequalities," *Studies in Computational Mathematics*, vol. 8, pp. 409–422, 2001.
- [33] K. You and L. Xie, "Network topology and communication data rate for consensusability of discrete-time multi-agent systems," *IEEE Transactions on Automatic Control*, vol. 56, no. 10, pp. 2262–2275, 2011.
- [34] D. Varagnolo, F. Zanella, A. Cenedese, G. Pillonetto, and L. Schenato, "Newton-Raphson Consensus for Distributed Convex Optimization," *IEEE Transactions on Automatic Control*, vol. 61, no. 4, pp. 994–009, 2016.
- [35] K. Margellos, A. Falsone, S. Garatti, and M. Prandini, "Distributed constrained optimization and consensus in uncertain networks via proximal minimization," *arXiv preprint arXiv:1603.02239*, 2016.
- [36] Y. Lou, G. Shi, K. H. Johansson, and Y. Hong, "Approximate projected consensus for convex intersection computation: Convergence analysis and critical error angle," *IEEE Transactions on Automatic Control*, vol. 59, no. 7, pp. 1722–1736, 2014.
- [37] D. P. Bertsekas, *Convex Optimization Algorithms*. Athena Scientific, 2015.
- [38] T.-H. Chang, A. Nedich, and A. Scaglione, "Distributed constrained optimization by consensus-based primal-dual perturbation method," *IEEE Transactions on Automatic Control*, vol. 59, no. 6, pp. 1524–1538, 2014.
- [39] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin, "On the linear convergence of the ADMM in decentralized consensus optimization," *IEEE Transactions on Signal Processing*, vol. 62, no. 7, pp. 1750–1761, 2014.
- [40] T.-H. Chang, M. Hong, and X. Wang, "Multi-agent distributed optimization via inexact consensus ADMM," *IEEE Transactions on Signal Processing*, vol. 63, no. 2, pp. 482–497, 2015.
- [41] F. Iutzeler, P. Bianchi, P. Ciblat, and W. Hachem, "Explicit convergence rate of a distributed alternating direction method of multipliers," *IEEE Transactions on Automatic Control*, vol. 61, no. 4, pp. 892–904, 2016.
- [42] H. J. Kushner and G. G. Yin, *Stochastic Approximation and Recursive Algorithms and Applications*. Springer, 2003.
- [43] R. Horn and C. Johnson, *Matrix Analysis*. Cambridge University Press, 1985.
- [44] H. Robbins and D. Siegmund, "A convergence theorem for non negative almost supermartingales and some applications," in *Herbert Robbins Selected Papers*. Springer, 1985, pp. 111–135.
- [45] S. Boyd. (2017) Subgradient methods. [Online]. Available: https://stanford.edu/class/ee364b/lectures/subgrad_method_slides.pdf
- [46] L. El Ghaoui and H. Lebret, "Robust solutions to least-squares problems with uncertain data," *SIAM Journal on Matrix Analysis and Applications*, vol. 18, no. 4, pp. 1035–1064, 1997.
- [47] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, Belmont, Massachusetts, US, 2003.



Keyou You received the B.S. degree in Statistical Science from Sun Yat-sen University, Guangzhou, China, in 2007 and the Ph.D. degree in Electrical and Electronic Engineering from Nanyang Technological University (NTU), Singapore, in 2012. After briefly working as a Research Fellow at NTU, he joined Tsinghua University in Beijing, China where he is now an Associate Professor in the Department of Automation. He held visiting positions at Politecnico di Torino, The Hong Kong University of Science and Technology, The University of Melbourne and etc.

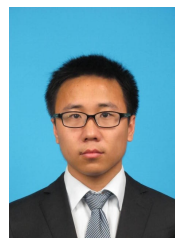
His current research interests include networked control systems, distributed algorithms, and their applications.

Dr. You received the Guan Zhaozhi award at the 29th Chinese Control Conference in 2010, and a CSC-IBM China Faculty Award in 2014. He was selected to the National 1000-Youth Talent Program of China in 2014, and received the National Science Fund for Excellent Young Scholars in 2017.



Roberto Tempo was a Director of Research of Systems and Computer Engineering at CNR-IEIIT, Politecnico di Torino, Italy. He held visiting positions at Tsinghua University in Beijing, Chinese Academy of Sciences, Kyoto University, The University of Tokyo, University of Illinois at Urbana-Champaign, German Aerospace Research Organization in Oberpfaffenhofen and Columbia University in New York. His research activities were focused on the analysis and design of complex systems with uncertainty, and various applications within information technology.

Dr. Tempo was a Fellow of the IEEE and a Fellow of the IFAC, a recipient of the IEEE Control Systems Magazine Outstanding Paper Award, of the Automatica Outstanding Paper Prize Award, and of the Distinguished Member Award from the IEEE Control Systems Society. He was a Corresponding Member of the Academy of Sciences, Institute of Bologna, Italy, Class Engineering Sciences. In 2010 Dr. Tempo was President of the IEEE Control Systems Society. He served as Editor-in-Chief of Automatica, Editor for Technical Notes and Correspondence of the IEEE Transactions on Automatic Control and Senior Editor of the same journal.



Pei Xie received the B.E. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2013. Currently, he is working toward his Ph.D. degree of the same institute. His research interests include networked control system, distributed optimizations, and their applications.