

A Stochastic Benders Decomposition Scheme for Large-Scale Data-Driven Network Design

Dimitris Bertsimas

Sloan School of Management and Operations Research Center, Massachusetts Institute of Technology, Cambridge, MA, US,
ORCID: 0000-0002-1985-1003
dbertsim@mit.edu

Ryan Cory-Wright

Department of Analytics, Marketing and Operations, Imperial College Business School, London, UK
IBM Thomas J. Watson Research Center, USA
ORCID: 0000-0002-4485-0619
r.cory-wright@imperial.ac.uk

Jean Pauphilet

Management Science and Operations, London Business School, London, UK
ORCID: 0000-0001-6352-0984
jpauphilet@london.edu

Periklis Petridis

Operations Research Center, Massachusetts Institute of Technology, Cambridge, MA, USA
ORCID: 0000-0002-1019-0763
periklis@mit.edu

Network design problems involve constructing edges in a transportation or supply chain network to minimize construction and daily operational costs. We study a data-driven version where operational costs are uncertain and estimated on historical data. This problem is computationally challenging, and instances with as few as 50 nodes cannot be solved to optimality by current decomposition techniques. We propose a stochastic variant of Benders decomposition that mitigates the high computational cost of generating each cut by sampling a subset of the data at each iteration and nonetheless generates deterministically valid cuts (as opposed to the probabilistically valid cuts frequently proposed in the stochastic optimization literature) via a dual averaging technique. We implement both single-cut and multi-cut variants of this Benders decomposition, as well as a k -cut variant that uses clustering of the historical scenarios. On instances with 100–200 nodes, our algorithm achieves 4-5% optimality gaps, compared with 13-16% for deterministic Benders schemes, and scales to instances with 700 nodes and 50 commodities within hours. Beyond network design, our strategy could be adapted to generic two-stage stochastic mixed-integer optimization problems where second-stage costs are estimated via a sample average.

Key words: Generalized Benders Decomposition; Network Design; Stochastic Integer Optimization

History:

1. Introduction

Network design is one of the most famous and frequently studied problems in the Operations Research literature, with widespread applications in logistics, air transportation (Barnhart et al. 2003), supply chains (Santoso et al. 2005, Pishvaei et al. 2014), telecommunications (Balakrishnan et al. 1991) and energy markets (Binato et al. 2001) among other domains. These problems are large-scale and involve uncertain parameters which reflect deviations between the forecast and realized utilization of a network, e.g., uncertain consumer demand in an air traffic control problem or uncertain renewable generation output in a capacity expansion problem. Moreover, we often have data on past realizations of the uncertain parameters. Unfortunately, despite the rapid advances in the scalability of branch-and-bound solvers over the past 25 years, data-driven network design problems with as few as 50 nodes are, to our knowledge, currently regarded as intractable and instead are solved via domain-specific approximation algorithms or heuristics (Crainic et al. 2021a).

To scale to network design problems with up to 50 nodes, the mixed-integer optimization (MIO) community has developed a suite of algorithms for mixed-integer nonlinear problems over the past 25 years, originating with the works of Ceria and Soares (1999), Stubbs and Mehrotra (1999) and refined by Frangioni and Gentile (2006, 2009), Günlük and Linderoth (2009) among others. Essentially, these methods tackle mixed-integer problems with logical constraints and a partially separable objective function, and enforce logical constraints implicitly via perspective functions, thus tightening the Boolean relaxation. Indeed, mixed-integer decomposition schemes that exploit these perspective reformulations often solve problems to optimality at sizes an order of magnitude larger than was previously possible; see Fischetti et al. (2017), Bertsimas et al. (2021) for related generalized Benders decomposition schemes.

In a different direction, the machine-learning community has enjoyed considerable success over the past 25 years in improving the scalability of unconstrained data-driven optimization. A common meta-approach is to modify a classical optimization algorithm to sample from an observed dataset at each iteration of the algorithm, and not consider the entire dataset as part of each iterate. Remarkably, each sample often conveys the same essential information as the entire dataset but can be processed multiple orders of magnitude faster. This sampling approach routinely produces a multiple-order-of-magnitude scalability improvement on classical optimization algorithms. Stochastic variants of first-order methods such as Stochastic Gradient Descent (SGD, Davis et al. 2020), the Stochastic Average Gradient method (Schmidt et al. 2017), or Adam (Kingma and Ba 2014) are currently considered to be state-of-the-art for unconstrained problems.

In this paper, we propose to embed a sampling technique within a Benders decomposition (Geoffrion 1972) scheme run on the perspective reformulation (Günlük and Linderoth 2009) of a network design problem. We demonstrate that this approach obtains bound gaps of 4–5% on instances with

100–200 nodes, three times smaller than the bound gaps obtained by deterministic Benders decomposition schemes in a comparable amount of time. Moreover, our approach successfully scales to obtain bound gaps of 40% on instances with 700 nodes and 50 commodities. At this scale, deterministic Benders schemes are too expensive to obtain feasible solutions. Our numerical success can be explained by the fact that sampling allows us to generate significantly more Benders cuts within a given time budget than is possible via a deterministic Benders approach, while conveying most of the essential information stored in each deterministic cut. Although developed for the special case of data-driven multi-commodity capacitated fixed-charge network design problems, we believe our approach could be applied to two-stage stochastic optimization problems where the first-stage variables are discrete, and the second-stage cost is evaluated via a sample average approximation.

1.1. Problem Formulation and Main Contributions

We propose a new approach for solving data-driven Multi-commodity Capacitated Fixed-charge Network Design (MCFND) problems to certifiable optimality. Similar models appear in Magnanti and Wong (1984), Costa (2005), Rahmaniani et al. (2018) among other works.

In MCFND problems, there is an index set of commodities \mathcal{K} to be shipped over a capacitated directed network $(\mathcal{N}, \mathcal{E})$, where \mathcal{N} denotes a set of nodes and \mathcal{E} denotes a set of edges. Our overall objective is to perform this transshipment in a manner that minimizes the construction plus flow transportation cost. Let \mathbf{A} denote this network's corresponding flow conservation matrix. The capacity of arc $(i, j) \in \mathcal{E}$ is given by $u_{i,j}$ and each node $n \in \mathcal{N}$ supplies or demands an amount $d_n^{k,r}$ of each commodity $k \in \mathcal{K}$ in each scenario $r \in \mathcal{R}$. There is a fixed cost c_{ij} of activating each edge $(i, j) \in \mathcal{E}$, and given this problem data, we introduce binary design variables $z_{i,j} \in \{0, 1\}$ to denote whether the (i, j) th edge is activated. The flow variable $x_{ij}^{k,r}$ then denotes the quantity of commodity k routed on edge (i, j) in scenario r , and f_{ij}^k denotes the marginal transportation cost, i.e., the per unit cost of transporting the k th commodity through edge (i, j) . Moreover, we follow the standard sample-average-approximation paradigm (see Shapiro et al. 2021, for a general theory) in placing equal weight on each observation of historical data r in our objective.

The complete optimization formulation for data-driven MCFND can then be written as:

$$\begin{aligned}
\min \quad & \sum_{(i,j) \in \mathcal{E}} c_{i,j} z_{i,j} + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \sum_{(i,j) \in \mathcal{E}} \left(\sum_{k \in \mathcal{K}} f_{i,j}^k x_{i,j}^{k,r} + \frac{1}{2\gamma} \left(\sum_{k \in \mathcal{K}} x_{i,j}^{k,r} \right)^2 \right) \\
\text{s.t.} \quad & \mathbf{A} \mathbf{x}^{k,r} = \mathbf{d}^{k,r}, \quad \forall k \in \mathcal{K}, r \in \mathcal{R}, \\
& \sum_{k \in \mathcal{K}} x_{i,j}^{k,r} \leq u_{i,j}, \quad \forall (i,j) \in \mathcal{E}, r \in \mathcal{R}, \\
& \mathbf{x}^{k,r} \geq 0, \quad x_{i,j}^{k,r} = 0 \text{ if } z_{i,j} = 0, \quad \forall (i,j) \in \mathcal{E}, \\
& \sum_{(i,j) \in \mathcal{E}} z_{i,j} \leq c_0, \quad z_{i,j} \in \{0, 1\} \quad \forall (i,j) \in \mathcal{E},
\end{aligned} \tag{1}$$

where, for any $\gamma > 0$, $1/2\gamma \sum_{r \in \mathcal{R}} \left(\sum_{k \in \mathcal{K}} x_{i,j}^{k,r} \right)^2$ is a ridge regularization term that improves Problem (1)'s practical tractability (c.f. Bertsimas et al. 2021), by augmenting the hard constraint $\sum_{k \in \mathcal{K}} x_{i,j}^{k,r} \leq u_{i,j}$ on each edge's capacity with a soft penalty in the objective (see also Atamtürk and Günlük (2018) for a discussion of capacity constraints in network design problems). Correspondingly, Problem (1) is an upper approximation of the optimal objective value without regularization. Therefore, one would ideally like to minimize (1) with $1/\gamma \rightarrow 0$. However, since (1) is a data-driven formulation, the regularization term may be beneficial from a model stability perspective, particularly in settings with limited or noisy data.

In this paper, we provide two main contributions. First, we propose a new decomposition method that combines sampling-based methods from the stochastic optimization and machine learning literature with a Generalized Benders Decomposition approach in the spirit of Geoffrion (1972). Our approach is capable of tackling large-scale mixed-integer problems by leveraging weak duality to obtain valid dual variables for scenarios we do not explicitly sample.

Second, we implement and benchmark our approach across a wide variety of large-scale network design instances, and explore the performance benefits of various design and implementation choices. Our approach allows us to solve network design problems with 200 nodes to within 4-5% of optimality in hours, and obtain high-quality feasible solutions on instances with up to 700 nodes.

1.2. Background and Literature Review

Our work is built on two intertwined literatures. First, decomposition schemes for large-scale deterministic problems with logical constraints developed by the MIO community. Second, sampling algorithms for problems with exogenous uncertainty developed by the stochastic optimization community. We further remark that, owing to Problem (1)'s significant computational difficulty, a wide variety of approximation algorithms (Agrawal et al. 1991, Goemans and Bertsimas 1993, Bertsimas and Teo 1998) and heuristic methods have also been proposed for solving Problem (1); see Rodríguez-Martín and Salazar-González (2010), Gendron et al. (2018) for reviews.

Cutting-Plane Schemes for Mixed-Integer Optimization: Problem (1) is a computationally challenging mixed-integer problem that encompasses hard combinatorial problems such as Steiner tree optimization (Garey and Johnson 1977) and possesses extremely poor Boolean relaxations (Gendron et al. 1999). Indeed, generic branch-and-bound solvers cannot currently solve network design (ND) problems at even moderate problem sizes with tens of nodes (see Crainic et al. 2021b, Section 6.1 for an investigation of CPLEX version 12.8's performance on synthetic ND instances with ten nodes). Accordingly, and due to its cardinal importance in practice, ND has emerged as one of the most frequently studied problems in the MIO literature over the past 50 years.

Throughout the first 30 years of the field of Operations Research, there was a spirited debate regarding the most efficient technique for solving ND problems, with many proposals, including branch-and-bound (Boyce et al. 1973), Lagrangian methods (Cornuejols et al. 1980), and dynamic programming (Erickson et al. 1987). The idea of solving ND problems via Generalized Benders decomposition (Geoffrion 1972) was moved front-and-center by Magnanti and Wong (1981, 1984). Building upon several influential prior works, including Geoffrion and Graves (1974), Florian et al. (1976), Richardson (1976), they found that an accelerated Benders decomposition was a viable and often more scalable alternative for ND problems than several other optimization approaches, including the three aforementioned ones. Ever since, Benders decomposition has been widely recognized as one of the most competitive methods for solving ND problems; we refer to Fischetti et al. (2017), Crainic et al. (2021b) for modern reviews of Benders decomposition for ND problems.

In a related direction, a significant line of work has developed a suite of cutting planes that iteratively strengthen Problem (1)’s Boolean relaxation upon their imposition; see, e.g., Van Roy and Wolsey (1985), Magnanti et al. (1993, 1995), Bienstock et al. (1998), Günlük (1999), Atamtürk and Günlük (2021) and references therein. Remarkably, these approaches are so numerically successful and easy to implement that they are usually incorporated within commercial branch-and-cut solvers within several years of their proposal (Bixby 2012). As a result, some of the decomposition schemes reviewed above may even be considered “sleeping beauties” in the sense of Ke et al. (2015), i.e., were not originally considered numerically successful but would be if proposed today, implicitly in conjunction with these valid inequalities.

Decomposition Schemes for Large-Scale Optimization Under Uncertainty: Cotemporally, a considerable amount of attention has been devoted by the stochastic optimization community to solving large-scale convex optimization problems with uncertain parameters for which we have access to either a joint probability distribution or observations from historical data. Initiated by the independent works of Dantzig (1955), Beale (1955), and subsequently refined by Wets (1966), Van Slyke and Wets (1969), contemporary optimizers for large-scale stochastic problems typically invoke the Minkowski-Weyl theorem (c.f. Bertsimas and Tsitsiklis 1997, Chapter 4) to solve their deterministic equivalents via Benders decomposition (which was termed the L-shaped method by Van Slyke and Wets 1969). Alternatively, works like Zakeri et al. (2000), Fábíán (2000), Rei et al. (2009), Guigues (2020) propose generating Benders cuts without solving each subproblem to optimality.

The two main variants of Benders decomposition invoked for two-stage stochastic integer optimization problems such as Problem (1) are called single-cut and multi-cut Benders. Single-cut schemes maintain a single epigraph variable that upper bounds the expected transshipment cost and generate a single cut at each iteration of Benders decomposition. Multi-cut schemes associate a separate epigraph variable with the cost incurred in each scenario and generate a separate cut for

each epigraph variable in each iteration (Birge and Louveaux 1988). Therefore, single-cut schemes typically require more iterations to converge but require less time to perform each iteration (see Birge and Louveaux 2011, de Camargo et al. 2008, You and Grossmann 2013, for comparisons). Problems with fewer scenarios are typically solved faster via multi-cut approaches but the relative performance of each variant is highly problem-dependent.

More recently, a considerable amount of attention has been devoted to designing variants of Benders decomposition that avoid solving a subproblem for each scenario at each iteration by sampling. Hige and Sen (1991, 1996b), Pereira and Pinto (1991), Dantzig and Infanger (1993), Infanger (1992) initiated this line of inquiry by proposing stochastic cutting-plane schemes that converge almost surely (see also Bertsimas and Li 2022). Determining convergence of these schemes is technically challenging. Various statistical tests exist (see, e.g., Hige and Sen 1996a, Morton 1998, Mak et al. 1999) that provide confidence intervals on the duality gap. Yet, to avoid multiple-testing problems, practitioners typically run stochastic cutting-plane methods for a prespecified number of iterations and then perform a statistical test on termination (De Matos et al. 2015).

1.3. Structure

We propose a stochastic Benders decomposition scheme that combines the perspective reformulation technique from the MIO literature with sampling ideas from the stochastic optimization literature to, for the first time, successfully solve data-driven capacitated network design problems with hundreds of nodes to certifiable (near) optimality. The rest of this paper is laid out as follows:

- In Section 2, we propose stochastic variants of the single-, k -, and multi-cut versions of Benders decomposition to solve a perspective reformulation of (1). Our algorithms randomly sample a subset of scenarios $\mathcal{R}_t \subseteq \mathcal{R}$ at each iteration and use a dual averaging technique to generate cuts that are deterministically valid for all $r \in \mathcal{R}$, while previous stochastic approaches generate cuts that are only valid on average or with high probability. We prove high probability bounds on the approximation error stemming from our dual averaging technique.
- In Section 3, we propose rigorous convergence criteria to terminate our stochastic decomposition schemes at a certifiable optimal solution. Since our master optimization problem is a MIO problem, we also discuss the specific termination challenges arising when Benders Decomposition is implemented via branch-and-cut (or lazy constraints). We also propose techniques for accelerating the convergence of our methods, by warm-starting their upper and lower bounds.
- In Section 4, we apply our decomposition schemes to a collection of network design instances that are synthetically generated or obtained from the literature (Crainic et al. 2016, 2021a). On the largest instances solvable by a deterministic Benders decomposition algorithm, our stochastic cutting-plane strategy obtains optimality gaps three to four times smaller than their deterministic

counterparts. Moreover, our approach scales to instances with up to 700 nodes. On instances with around 200 nodes, it routinely achieves bound gaps on the order of 5% within two hours.

Notation

We let non-boldface characters such as b denote scalars, lowercase bold-faced characters (\mathbf{x}) denote vectors, uppercase bold-faced characters (\mathbf{A}) denote matrices, and calligraphic uppercase characters (\mathcal{Z}) denote sets. We let $[n]$ denote the running set of indices $\{1, \dots, n\}$. We let \mathbf{e} denote the vector of ones, and $\mathbf{0}$ denote the vector of all zeros. Finally, we let $\langle \cdot, \cdot \rangle$ denote the Euclidean inner product between two vectors of the same size, i.e., $\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{i=1}^n x_i y_i$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

2. Deterministic and Stochastic Cutting-Plane Methods

In this section, we propose an efficient numerical strategy for solving Problem (1) to certifiable optimality. The backbone of our approach is a generalized Benders decomposition scheme run on a perspective reformulation of Problem (1), which uses sampling techniques to avoid explicitly solving each scenario at each iteration of the method. Instead, we use dual-optimal solutions of the sampled subproblems to construct dual-feasible solutions to the remaining subproblems and thereby construct valid cuts. We further discuss the convergence properties of our method.

2.1. A Two-Stage Reformulation

We open this section by observing that the flow minimization problem with respect to each $\mathbf{x}^{:r}$ in (1) is decomposable across scenarios $r \in \mathcal{R}$. Therefore, consider a set of demand vectors $\mathbf{d}^k \in \mathbb{R}^{\mathcal{N}}$ for $k \in \mathcal{K}$ and define

$$f(\mathbf{z}; \mathbf{d}) := \min_{\mathbf{x}^k \in \mathbb{R}_+^{\mathcal{E}}, k \in \mathcal{K}} \sum_{k \in \mathcal{K}} \langle \mathbf{f}^k, \mathbf{x}^k \rangle + \frac{1}{2\gamma} \sum_{(i,j) \in \mathcal{E}} \left(\sum_{k \in \mathcal{K}} x_{i,j}^k \right)^2 \quad \text{s.t. } \mathbf{A}\mathbf{x}^k = \mathbf{d}^k, \forall k \in \mathcal{K}, \quad (2)$$

$$\sum_{k \in \mathcal{K}} x_{i,j}^k \leq u_{i,j}, \forall (i,j) \in \mathcal{E},$$

$$x_{i,j}^k = 0 \text{ if } z_{i,j} = 0, \forall (i,j) \in \mathcal{E},$$

to be the operational cost of serving demand \mathbf{d} on network $(\mathcal{N}, \mathcal{E})$ with design variables \mathbf{z} . Observe that the minimization problem defining $f(\mathbf{z}; \mathbf{d})$ is not decomposable across commodities because of shared capacity constraints. With this notation, Problem (1) is equivalent to

$$\min_{\mathbf{z} \in \mathcal{Z}} \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} f(\mathbf{z}; \mathbf{d}^{:r}), \quad (3)$$

where $\mathbf{d}^{:r}$ denotes the collection of demand vectors $\{\mathbf{d}^{k,r}, k \in \mathcal{K}\}$ and $\mathcal{Z} = \{\mathbf{z} \in \{0, 1\}^{\mathcal{E}} : \sum_{(i,j)} z_{i,j} \leq c_0\}$ denotes the set of feasible edges. The network design formulation (3) separates the discrete design variables \mathbf{z} from the continuous second-stage routing variables \mathbf{x}^k , thus giving a pure integer optimization formulation that is readily amenable to outer-approximation techniques.

2.2. A Linear Lower Approximation of the Second-Stage Cost Function

In this section, we derive a family of Benders cuts that successfully outer-approximate a perspective reformulation of (2). We remark that this derivation is now standard in the literature (see, e.g., Bertsimas et al. 2021), although we provide it here to keep this work self-contained.

Since the objective function in (3) involves the average of the function $f(\mathbf{z}, \mathbf{d})$ over $|\mathcal{R}|$ realization of \mathbf{d} , we start by analyzing properties of the function $f(\mathbf{z}, \mathbf{d})$ in isolation, with a view to establish that $f(\mathbf{z}, \mathbf{d})$ is convex in \mathbf{z} and a valid subgradient can be obtained by solving a dual problem.

PROPOSITION 1. *For any $\mathbf{z} \in \{0, 1\}^\mathcal{E}$ and demand vectors \mathbf{d}^k , $k \in \mathcal{K}$ such that Problem (2) admits a feasible solution, we have:*

$$f(\mathbf{z}; \mathbf{d}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^\mathcal{E}, \boldsymbol{\beta} \in \mathbb{R}_+^\mathcal{E} \\ \mathbf{p}^k \in \mathbb{R}^\mathcal{N}, k \in \mathcal{K}}} \sum_{k \in \mathcal{K}} \langle \mathbf{p}^k, \mathbf{d}^k \rangle - \langle \boldsymbol{\beta}, \mathbf{u} \rangle - \frac{\gamma}{2} \sum_{(i,j) \in \mathcal{E}} z_{i,j} (\alpha_{i,j} + \beta_{i,j})^2 \text{ s.t. } \mathbf{A}^\top \mathbf{p}^k \leq \mathbf{f}^k - \boldsymbol{\alpha}. \quad (4)$$

The proof of Proposition 1 follows analogously to Bertsimas et al. (2021, Theorem 2.5) and relies on deriving the dual of the minimization problem defining $f(\mathbf{z}; \mathbf{d})$ by using a variable decomposition *à la Fenchel*; for completeness, we provide a formal proof in Appendix A.1.

Proposition 1 calls for a few observations. First, according to the dual reformulation, $f(\mathbf{z}; \mathbf{d})$ can be expressed as the point-wise maximum of affine functions in \mathbf{z} , hence $f(\mathbf{z}; \mathbf{d})$ is convex in \mathbf{z} . Second, any feasible dual solution $\boldsymbol{\alpha} \in \mathbb{R}^\mathcal{E}$, $\boldsymbol{\beta} \in \mathbb{R}_+^\mathcal{E}$, $\mathbf{p}^k \in \mathbb{R}^\mathcal{N}$ such that $\mathbf{A}^\top \mathbf{p}^k \leq \mathbf{f}^k - \boldsymbol{\alpha}$ provides a valid linear lower approximation of $f(\mathbf{z}'; \mathbf{d})$. Namely, for any \mathbf{z}' ,

$$f(\mathbf{z}'; \mathbf{d}) \geq \sum_{k \in \mathcal{K}} \langle \mathbf{p}^k, \mathbf{d}^k \rangle - \langle \boldsymbol{\beta}, \mathbf{u} \rangle - \frac{\gamma}{2} \sum_{(i,j) \in \mathcal{E}} z'_{i,j} (\alpha_{i,j} + \beta_{i,j})^2.$$

When the dual variables are optimal for a particular vector \mathbf{z} , the resulting offset and slope in the above linear approximation are exactly the value of $f(\mathbf{z}; \mathbf{d})$ and a subgradient of f at \mathbf{z} , i.e.,

$$f(\mathbf{z}'; \mathbf{d}) \geq f(\mathbf{z}; \mathbf{d}) + \langle \nabla f(\mathbf{z}; \mathbf{d}), \mathbf{z} - \mathbf{z}' \rangle.$$

Third, Proposition 1 applies if Problem (2) is feasible for the current design vector \mathbf{z} . On the other hand, if (2) is not feasible, then the following feasibility problem does not admit a solution:

$$\exists \mathbf{x} \in \mathbb{R}_+^{\mathcal{E} \times \mathcal{K}} : \mathbf{A}\mathbf{x}^k = \mathbf{d}^k \quad \forall k \in \mathcal{K}, \quad \sum_{k \in \mathcal{K}} x_{i,j}^k \leq u_{i,j} z_{i,j}, \quad \forall (i,j) \in \mathcal{E}.$$

Hence, by Farkas's lemma (see, e.g., Bertsimas and Tsitsiklis 1997, Theorem 4.6), we can find a certificate of infeasibility, i.e., $\boldsymbol{\beta} \in \mathbb{R}_+^\mathcal{E}$, $\mathbf{p}^k \in \mathbb{R}^\mathcal{N}$, $k \in \mathcal{K}$ such that $\langle \mathbf{p}^k, \mathbf{d}^k \rangle > \sum_{(i,j) \in \mathcal{E}} z_{i,j} u_{i,j} \beta_{i,j}$. Therefore, we can separate \mathbf{z} from the set of feasible design vectors \mathbf{z} by imposing the feasibility cut

$$\sum_{k \in \mathcal{K}} \langle \mathbf{p}^k, \mathbf{d}^k \rangle \leq \sum_{(i,j) \in \mathcal{E}} z'_{i,j} u_{i,j} \beta_{i,j}. \quad (5)$$

Finally, as has already been observed in the literature (Xie and Deng 2020, Bertsimas et al. 2021), our reformulation can alternatively be achieved by performing a perspective reformulation on (2) to rewrite it as a mixed-integer second-order cone problem (c.f. Günlük and Linderoth 2009) and taking the dual of this perspective reformulation with respect to the continuous variables.

2.3. Epigraph Formulations: Modeling Choice and Algorithmic Implications

In this section, we exploit our previously developed characterization of $f(\mathbf{z}, \mathbf{d})$ as the pointwise maximum of functions linear in \mathbf{z} to revisit three deterministic outer-approximation methods that solve Problem (1) to certifiable optimality. For simplicity, we only treat the optimality cuts in this section; feasibility cuts follow in much the same way.

Outer-approximation methods such as generalized Benders decomposition solve (3) by constructing a lower approximation of the second-stage operational cost $\frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} f(\mathbf{z}; \mathbf{d}^{:r})$ and refining this approximation at each step. However, since the second-stage cost is the average operational cost over $|\mathcal{R}|$ scenarios, one can either approximate each term $f(\mathbf{z}; \mathbf{d}^{:r})$ separately, their sum, or k partial sums separately, which we refer to as multi-cut, single-cut, and k -cut approaches, respectively.

In a multi-cut approach, we consider the following epigraph formulation of Problem (3), as originally proposed by Birge and Louveaux (1988) for two-stage stochastic linear optimization:

$$\min_{\substack{\mathbf{z} \in \mathcal{Z} \\ \eta_r \in \mathbb{R}, \forall r \in \mathcal{R}}} \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \eta_r \text{ s.t. } \eta_r \geq f(\mathbf{z}; \mathbf{d}^{:r}), \forall r \in \mathcal{R},$$

and iteratively refine a piecewise linear lower approximation of $f(\mathbf{z}; \mathbf{d}^{:r})$ for each epigraph constraints until convergence. Specifically, at each iteration T , the multi-cut cutting-plane algorithm solves the MIO problem

$$\min_{\substack{\mathbf{z} \in \mathcal{Z} \\ \eta_r \in \mathbb{R}, \forall r \in \mathcal{R}}} \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \eta_r \text{ s.t. } \eta_r \geq f(\mathbf{z}^t; \mathbf{d}^{:r}) + \langle \nabla f(\mathbf{z}^t; \mathbf{d}^{:r}), \mathbf{z} - \mathbf{z}^t \rangle, \forall t \in [T], \forall r \in \mathcal{R}. \quad (6)$$

Observe that, in this implementation, each of the $|\mathcal{R}|$ functions $f(\mathbf{z}; \mathbf{d}^{:r})$ is linearized at T points \mathbf{z}^t , so (6) comprises $|\mathcal{R}| \times T$ linear constraints. The solution of (6), \mathbf{z}^{T+1} , then serves as a linearization point to further improve the approximations of the functions $f(\mathbf{z}; \mathbf{d}^{:r})$ at the next iteration.

Alternatively, the single-cut approach, as originally proposed for two-stage stochastic linear optimization by Van Slyke and Wets (1969), considers a more compact epigraph formulation:

$$\min_{\substack{\mathbf{z} \in \mathcal{Z} \\ \eta \in \mathbb{R}}} \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \eta \text{ s.t. } \eta \geq \sum_{r \in \mathcal{R}} f(\mathbf{z}; \mathbf{d}^{:r}),$$

and constructs a piece-wise linear lower-approximation of $\sum_{r \in \mathcal{R}} f(\mathbf{z}; \mathbf{d}^{:r})$ directly. In a single-cut cutting-plane algorithm, at a given iteration T , the epigraph constraint is replaced by the linear constraints of the form

$$\eta \geq \sum_{r \in \mathcal{R}} f(\mathbf{z}^t; \mathbf{d}^{:r}) + \left\langle \sum_{r \in \mathcal{R}} \nabla f(\mathbf{z}^t; \mathbf{d}^{:r}), \mathbf{z} - \mathbf{z}^t \right\rangle. \quad (7)$$

The single-cut approach involves only one epigraph variable η (compared with $|\mathcal{R}|$ in the multi-cut implementation) and adds one linear constraint at each iteration (vs. $|\mathcal{R}|$). As a result, the MIO problems involved in the single-cut approach are smaller and usually more tractable than those solved by the multi-cut approach. Yet, multi-cut methods approximate the second-stage cost function more accurately and therefore might require fewer iterations to converge. Various studies, including Birge and Louveaux (2011), de Camargo et al. (2008), You and Grossmann (2013) have reported mixed results on the relative merits of single and multi-cut methods, and which method works best appears to depend on the underlying problem and the number of scenarios.

To successfully combine the best aspects of single and multi-cut approaches, a k -cut approach was proposed by Trukhanov et al. (2010), Contreras et al. (2011). They observed that scenarios can often be partitioned into subsets (or clusters) that are very similar to one another. Moreover, aggregating the cuts in each partition successfully compresses information about the second-stage cost surface and, on a per-iteration basis, is almost as fast as a single-cut approach. Accordingly, let $\cup_{c \in [k]} \mathcal{S}_c$ be a partition of \mathcal{R} . Then, at each iteration, the k -cut approach solves the MIO:

$$\min_{\substack{\mathbf{z} \in \mathcal{Z} \\ \eta_c \in \mathbb{R}, c \in [k]}} \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \sum_{c \in [k]} \eta_c \text{ s.t. } \eta_c \geq \sum_{r \in \mathcal{S}_c} f(\mathbf{z}^t; \mathbf{d}^r) + \langle \nabla f(\mathbf{z}^t; \mathbf{d}^r), \mathbf{z} - \mathbf{z}^t \rangle, \forall t \in [T], \forall c \in [k], \quad (8)$$

and constructs each cut similarly to the single and multi-cut approaches. At each iteration, the k -cut approach adds k linear constraints (one per cluster $c \in [k]$). If $k = 1$ (resp. $|\mathcal{R}|$), we recover the single-cut (resp. multi-cut) algorithm. We remark that all three methods converge in a finite but possibly exponential number of iterations by the finiteness of $\{0, 1\}^\mathcal{E}$ and since no method visits a binary vector \mathbf{z} twice (see also Geoffrion 1972, Theorem 2.4).

A common thread between all three approaches is that evaluating values of functions of the form $f(\mathbf{z}, \mathbf{d})$ (and their subgradients) is the main computational bottleneck, and, in all three approaches, the number of function evaluations is the same, $|\mathcal{R}|$, which can be prohibitive, especially when the number of past scenarios $|\mathcal{R}|$ increases. Accordingly, we propose stochastic versions of these approaches with improved per-iteration computational complexity in the next section.

2.4. A Sample-Based Stochastic Cutting-Plane Algorithm

In this section, we propose stochastic variants of the cutting-plane methods proposed in the previous section, which obtain high-quality deterministically valid lower bounds without explicitly solving an optimization problem in each scenario and each commodity at each iteration of the method. We also discuss the convergence of these methods. As these methods do not provide deterministically valid upper bounds from a single sample, we defer a detailed discussion of their upper bounds, the corresponding termination criteria, and their single-tree implementation to Section 3, and assume for ease of exposition that all cutting-plane methods are multi-tree throughout the section.

instead of (7), where the quantities $\frac{|\mathcal{R}_t|}{|\mathcal{R}|} \sum_{r \in \mathcal{R}_t} f(\mathbf{z}^t; \mathbf{d}^{:r})$ and $\frac{|\mathcal{R}_t|}{|\mathcal{R}|} \sum_{r \in \mathcal{R}_t} \nabla f(\mathbf{z}^t; \mathbf{d}^{:r})$ are unbiased estimates of the original offset and slope terms, $\sum_{r \in \mathcal{R}} f(\mathbf{z}^t; \mathbf{d}^{:r})$ and $\sum_{r \in \mathcal{R}} \nabla f(\mathbf{z}^t; \mathbf{d}^{:r})$ respectively, so that (9) is a reasonable approximation of the original constraint (7). This intuition is similar to that of SGD in unconstrained continuous optimization. Unfortunately, these cuts are only valid probabilistically and may actually cut off part of the feasible region when combined. Moreover, while the sampled cuts are unbiased estimates of the slope, optimizing these estimates via Benders decomposition yields solutions that disappoint significantly out-of-sample due to the so-called optimizer’s curse (Smith and Winkler 2006). SGD shares the same drawbacks but mitigates them by performing only one gradient step at each iteration and forgetting estimation errors between iterations. Conversely, in a cutting-plane algorithm, cuts added at one iteration are imposed in subsequent iterations, until termination.

We reconcile the computational benefits of sampling with the aforementioned drawbacks of the stochastic single-cut approach by leveraging the dual formulation of $f(\mathbf{z}; \mathbf{d})$ in Proposition 1 to derive deterministically valid lower-approximations for scenarios r that are not sampled. Further, we argue that provided the sampled scenarios are sufficiently representative of the remaining scenarios, this approximation is sufficiently accurate that we eventually obtain an optimal solution with high probability; see also Zakeri et al. (2000) for an “inexact” Benders decomposition method.

Specifically, recall that any feasible dual solution $(\boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{p})$ provides a valid lower bound:

$$f(\mathbf{z}; \mathbf{d}^{:r}) \geq q(\mathbf{z}^t, \boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{p}; \mathbf{d}^{:r}) + \langle \nabla_{\mathbf{z}} q(\mathbf{z}^t, \boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{p}; \mathbf{d}^{:r}), \mathbf{z} - \mathbf{z}^t \rangle,$$

with $q(\mathbf{z}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{p}; \mathbf{d}) := \sum_{k \in \mathcal{K}} \langle \mathbf{p}^k, \mathbf{d}^k \rangle - \langle \boldsymbol{\beta}, \mathbf{u} \rangle - \frac{\gamma}{2} \sum_{(i,j) \in \mathcal{E}} z_{i,j} (\alpha_{i,j} + \beta_{i,j})^2$. Hence, we replace (7) by a constraint of the form

$$\eta \geq \sum_{r \in \mathcal{R}} q(\mathbf{z}^t, \boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r; \mathbf{d}^{:r}) + \sum_{r \in \mathcal{R}} \langle \nabla_{\mathbf{z}} q(\mathbf{z}^t, \boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r; \mathbf{d}^{:r}), \mathbf{z} - \mathbf{z}^t \rangle, \quad (10)$$

for some feasible dual solutions $(\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r)$. Observe that, unlike (9), the constraint (10) is a deterministically valid (although not necessarily tight) lower bound on the true operational cost.

Collecting these observations yields our overall stochastic single-cut approach: First, to reduce the computational burden of solving an optimization problem for each scenario, at each iteration, we only solve a random subset of scenarios $r \in \mathcal{R}_t \subseteq \mathcal{R}$ —hence effectively computing $f(\mathbf{z}^t; \mathbf{d}^{:r})$ and $\nabla f(\mathbf{z}^t; \mathbf{d}^{:r})$. Second, for the remaining scenarios $r \notin \mathcal{R}_t$, we refrain from solving (4) and instead use the cheap to compute and feasible dual average solution $(\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r) = (\bar{\boldsymbol{\alpha}}^{\mathcal{R}_t}, \bar{\boldsymbol{\beta}}^{\mathcal{R}_t}, \bar{\mathbf{p}}^{\mathcal{R}_t}) := \frac{1}{|\mathcal{R}_t|} \sum_{r' \in \mathcal{R}_t} (\boldsymbol{\alpha}^{r'}, \boldsymbol{\beta}^{r'}, \mathbf{p}^{r'})$ instead. This gives a stochastic cutting-plane method with a sequence of deterministically valid non-decreasing lower bounds, which we formalize in Algorithm 2 (we defer a detailed discussion of its single-tree implementation and termination criterion to Section 3).

However, it is not obvious whether this method converges towards an optimal solution (e.g., in a limit) or whether it generates a never-ending sequence of deterministically valid but not tight cuts. We now provide some reassurance in this direction by showing that for the incumbent solution \mathbf{z}^t , the approximation error of cuts obtained via dual averaging can be decomposed, with high probability, as the sum of two terms: one term that depends on the variance of the optimal dual variables and that captures the heterogeneity in the demand scenarios, and one estimation error term that vanishes as $|\mathcal{R}_t|$ grows (proof deferred to Appendix A.2):

PROPOSITION 2. *Fix \mathbf{z}^t . For any $r \in \mathcal{R}$, denote $(\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r)$ the optimal dual solutions of (4) for $\mathbf{z} = \mathbf{z}^t$ and $\mathbf{d} = \mathbf{d}^{\cdot r}$. Denote ν^2 the variance in optimal dual variables, defined as*

$$\nu^2 = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left\| (\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r) - (\bar{\boldsymbol{\alpha}}^{\mathcal{R}}, \bar{\boldsymbol{\beta}}^{\mathcal{R}}, \bar{\mathbf{p}}^{\mathcal{R}}) \right\|^2 \quad \text{with } (\bar{\boldsymbol{\alpha}}^{\mathcal{R}}, \bar{\boldsymbol{\beta}}^{\mathcal{R}}, \bar{\mathbf{p}}^{\mathcal{R}}) = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} (\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r).$$

Then, there exist universal constants $L, M > 0$ such that, for any $\delta \in (0, e^{-1})$, when \mathcal{R}_t is sampled without replacement from \mathcal{R} with a fixed size $|\mathcal{R}_t|$, we have with probability $1 - 3\delta$:

$$\sum_{r \notin \mathcal{R}_t} \left| q(\mathbf{z}^t, \bar{\boldsymbol{\alpha}}^{\mathcal{R}_t}, \bar{\boldsymbol{\beta}}^{\mathcal{R}_t}, \bar{\mathbf{p}}^{\mathcal{R}_t}; \mathbf{d}^r) - f(\mathbf{z}^t; \mathbf{d}^r) \right| \leq L \sqrt{|\mathcal{R} \setminus \mathcal{R}_t|} \nu + D \sqrt{|\mathcal{R} \setminus \mathcal{R}_t| \log(1/\delta)}, \quad (11)$$

with

$$D := LM \sqrt{2|\mathcal{E}| + |\mathcal{N}| \times |\mathcal{K}|} \left[\sqrt{|\mathcal{R}|} \left(\frac{1}{|\mathcal{R}_t|} - \frac{1}{|\mathcal{R}|} \right)^{1/2} + \left(\frac{1}{|\mathcal{R} \setminus \mathcal{R}_t|} - \frac{1}{|\mathcal{R}|} \right)^{1/4} \right].$$

Algorithm 2 A Single-Cut Sample-Based Cutting Plane Method

- 1: **initialize** $\mathbf{z}_1; f(\mathbf{z}_0; \mathbf{d}^{\cdot r}), \nabla f(\mathbf{z}_0; \mathbf{d}^{\cdot r}), \forall r \in \mathcal{R}_0$.
 - 2: **set** $T \leftarrow 1$
 - 3: **repeat**
 - 4: **compute** $\mathbf{z}^{T+1}, \eta^{T+1} \leftarrow \arg \min_{\mathbf{z}, \eta} \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \eta$
 - s.t. $\eta \geq \sum_{r \in \mathcal{R}} q(\mathbf{z}^t; \mathbf{d}^{\cdot r}) + \langle \nabla q(\mathbf{z}^t; \mathbf{d}^{\cdot r}), \mathbf{z} - \mathbf{z}^t \rangle, \forall t \in [T]$,
 - 5: **sample** $\mathcal{R}_{T+1} \subseteq \mathcal{R}$
 - 6: **calculate** $f(\mathbf{z}^{T+1}; \mathbf{d}^{\cdot r}), \nabla f(\mathbf{z}^{T+1}; \mathbf{d}^{\cdot r})$ for $r \in \mathcal{R}_{T+1}$
 - 7: **set** $T \leftarrow T + 1$
 - 8: **until** Termination Criterion Met
-

Proposition 2 provides a probabilistic guarantee on the quality of each cut in terms of the sample size $|\mathcal{R}_t|$. Observe that the approximation error is proportional to $\sqrt{|\mathcal{R} \setminus \mathcal{R}_t|}$, which means that the approximation error is zero in the limit where $\mathcal{R}_t \rightarrow \mathcal{R}$ (as expected) but which also means that the approximation error grows sub-linearly in the number of scenarios to approximate $|\mathcal{R} \setminus \mathcal{R}_t|$.

Moreover, by the probabilistic method (see, e.g., Grimmett and Stirzaker 2020), Proposition 2 reveals that, for any \mathbf{z}^t and sufficiently small δ , there exists some \mathcal{R}_t such that this guarantee holds

deterministically. Indeed, setting $\delta < (1 - (|\mathcal{R}|/|\mathcal{R}_t|)^{-1})/3$ reveals that, with the notations of Proposition 2, repeatedly sampling \mathcal{R}_t for a given \mathbf{z}^t eventually gives a cut which is an underestimator of $f(\mathbf{z}^t)$ by at most ρ , where

$$\rho := L\sqrt{|\mathcal{R} \setminus \mathcal{R}_t|}\nu + D\sqrt{|\mathcal{R} \setminus \mathcal{R}_t| \log(1/\delta)}. \quad (12)$$

The above observation implies that running Algorithm 2 without termination and selecting a \mathbf{z}^t , which minimizes our underestimator in the limit, almost surely returns a ρ -optimal solution to Problem (1), where ρ is defined by Equation (12). Therefore, in practice, when Algorithm 2's lower bound stabilizes, we can either increase the number of scenarios sampled (and thus reduce ρ), or terminate with confidence if, according to a statistical test, the gap between our stochastic upper bound (see Section 3) and our deterministic lower bound is sufficiently small. As we observe in our numerical results (see Section 4), the optimality gap from single-cut at termination with a sample rate of around 10% is usually quite small in practice.

Finally, a stochastic variant of the k -cut approach can be developed analogously, by applying our method for stochastic single-cut to each cluster $c \in [k]$. Namely, we partition the set of scenarios \mathcal{R} into k sets $\mathcal{S}_c : c \in [k]$, and impose valid constraints of the form (10) to each epigraph variable η_c :

$$\eta_c \geq \sum_{r \in \mathcal{S}_c} q(\mathbf{z}^t, \boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r; \mathbf{d}^{r'}) + \sum_{r \in \mathcal{S}_c} \langle \nabla_{\mathbf{z}} q(\mathbf{z}^t, \boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r; \mathbf{d}^{r'}), \mathbf{z} - \mathbf{z}^t \rangle, \quad \forall c \in [k]. \quad (13)$$

Then, at each iteration, we sample and solve (4) for a random subset $\mathcal{R}_{t,c} \subseteq \mathcal{S}_c$ of scenarios in each cluster and set $(\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r) = (\bar{\boldsymbol{\alpha}}^{\mathcal{R}_{r,c}}, \bar{\boldsymbol{\beta}}^{\mathcal{R}_{r,c}}, \bar{\mathbf{p}}^{\mathcal{R}_{r,c}})$ for $r \in \mathcal{S}_c \setminus \mathcal{R}_{t,c}$. From Proposition 2 applied to each cluster separately, we obtain that the approximation error for cluster c is bounded, with high probability, by a term that depends on the variance in dual optimal variables within cluster c , ν_c^2 , plus a bootstrap estimation error term. Hence, if the clustering successfully reduces total weighted variance $\sum_{c \in [k]} \sqrt{|\mathcal{S}_c|} \nu_c$, a k -cut approach could improve the lower bound obtained by single-cut, while using the same number of samples per iteration.

In practice (and in our implementation), it is not feasible to cluster the set of scenarios \mathcal{R} based on their associated optimal dual solutions $\boldsymbol{\alpha}^r$, because the optimal dual solution implicitly depends on the current incumbent solution. Instead, we apply a k -means algorithm on the demand vectors $\mathbf{d}^{r'}$. The intuition is that the optimization problem defining $\boldsymbol{\alpha}^r$, (4), is parametrized by $\mathbf{d}^{r'}$. From sensitivity analysis, the optimal solution should be (irrespective of the incumbent solution) a smooth function of the demand vectors. So, low variance in terms of demand vectors $\mathbf{d}^{r'}$ should translate into low variance in terms of optimal solutions.

REMARK 1. Although the dual averaging technique is not needed to develop a stochastic multi-cut cutting-plane algorithm, it can be used to improve its convergence. In Algorithm 1, instead of only imposing a new cut for the epigraph variables η_r with $r \in \mathcal{R}_t$, we can also use dual averaging to impose new valid cuts on $\eta_r, r \notin \mathcal{R}_t$ as well.

3. Upper Bounds in Stochastic Cutting Planes with Binary Variables

In this section, we analyze the upper bounds obtained at each iteration of our cutting-plane methods and design convergence criteria that allow us to terminate our methods with confidence.

The primary motivation for this section is that while the lower bounds for the three stochastic cutting plane methods introduced in Section 2 are deterministic, their per iteration estimates of the cost associated with each incumbent solution \mathbf{z}^t ,

$$\langle \mathbf{c}, \mathbf{z}^t \rangle + \frac{1}{|\mathcal{R}_t|} \sum_{r \in \mathcal{R}_t} f(\mathbf{z}^t; \mathbf{d}^{:,r}) \quad (14)$$

are stochastic estimates that depend on the sample \mathcal{R}_t . Accordingly, we cannot simply use these stochastic estimates in the same way as in a deterministic method and terminate when the deterministic lower bound, say $\langle \mathbf{c}, \mathbf{z}^t \rangle + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \eta_r$ in the multi-cut case, is within ϵ of our stochastic upper bound, or we may terminate because \mathbf{z}^t is a high variance solution and we picked an optimistic sample set \mathcal{R}_t , rather than because \mathbf{z}^t is an optimal solution¹; see also Smith and Winkler (2006).

In addition, another salient characteristic of our problem is that the decision variables \mathbf{z} are binary. Hence, as described in pseudo-code in Algorithm 1 and 2, a MIO problem needs to be solved at each iteration by constructing a branch-and-bound tree (multi-tree implementation). Nowadays, efficient implementations of these schemes exist that simultaneously construct the branch-and-bound tree and generate cutting planes (single-tree implementation). We also discuss the extent to which the stochastic cutting-plane algorithms we developed in the previous section can be implemented with a single-tree instead of multi-tree approach.

3.1. Convergence Criteria

In this section, we define a convergence criterion by using an asymptotically normal estimator of the upper bound and using a related upper confidence bound. Suppose that one of our stochastic cutting-plane methods finds a solution \mathbf{z} , and that we would like to evaluate its quality. Then, we can use a sample \mathcal{W} to estimate the true cost of this solution

$$\bar{c} = \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} f(\mathbf{z}; \mathbf{d}^{:,r})$$

by its estimate on \mathcal{W} :

$$\hat{c}^{\mathcal{W}} = \langle \mathbf{c}, \mathbf{z} \rangle + \frac{1}{|\mathcal{W}|} \sum_{r \in \mathcal{W}} f(\mathbf{z}; \mathbf{d}^{:,r}).$$

In this section, for simplicity, we omit the dependency of $\hat{c}^{\mathcal{W}}$, \bar{c} , and the following quantities, on the solution \mathbf{z} . We also denote \mathcal{W} the random sample used for termination since it could be a new

¹ We note that such a deterministic convergence criterion is used in the stochastic method of Bertsimas and Li (2022).

independent draw from the sample \mathcal{R}_t used in the algorithm (and should be, for our estimation procedure to be unbiased).

As noted by Morton (1998), Mak et al. (1999), under some mild assumptions on the distribution of \mathbf{d}^k (e.g., finite variance), for an infinite number of scenarios $|\mathcal{R}|$, this estimator obeys a central limit theorem:

$$\sqrt{|\mathcal{W}|} [\hat{c}^{\mathcal{W}} - \bar{c}] \xrightarrow{d} \mathcal{N}(0, \sigma_c^2) \text{ as } |\mathcal{W}| \rightarrow \infty,$$

where $\sigma_c^2 = \text{Var}(f(\mathbf{z}, \mathbf{d}^r))$ can be estimated via the sample variance estimator

$$\hat{\sigma}_c^2 := \frac{1}{|\mathcal{W}| - 1} \sum_{r \in \mathcal{W}} \left(f(\mathbf{z}, \mathbf{d}^r) - \frac{1}{|\mathcal{W}|} \sum_{s \in \mathcal{W}} f(\mathbf{z}; \mathbf{d}^s) \right)^2.$$

In reality, however, we only have finitely many observations \mathcal{R} . Yet, provided $|\mathcal{R}|$ is large relative to $|\mathcal{W}|$, we can still apply the CLT to estimate the cost of \mathbf{z} . Consequently, letting q_α be such that $\mathbb{P}(\mathcal{N}(0, 1) \leq q_\alpha) = 1 - \alpha$, we can construct an asymptotically valid confidence interval for this estimator at level α of the form

$$\left[\hat{c}^{\mathcal{W}} - \frac{q_{\alpha/2}}{\sqrt{|\mathcal{W}|}} \hat{\sigma}_c, \hat{c}^{\mathcal{W}} + \frac{q_{\alpha/2}}{\sqrt{|\mathcal{W}|}} \hat{\sigma}_c \right].$$

We terminate our method using a modified version of the convergence criteria proposed by Morton (1998). Namely, letting

$$\bar{c}_{\alpha,t} := \hat{c}^{\mathcal{W}} + \frac{q_{\alpha/2}}{\sqrt{|\mathcal{W}|}} \hat{\sigma}_c$$

denote an upper confidence bound at level α on the cost of \mathbf{z}^t , the solution generated at the t th iterate of one of our cutting-plane methods, we terminate as soon the conservative bound gap falls below a predefined threshold ϵ , i.e., for the multi-cut method

$$\frac{\bar{c}_{\alpha,t} - \left(\langle \mathbf{c}, \mathbf{z}^t \rangle + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \eta_{r,t} \right)}{\bar{c}_{\alpha,t}} \leq \epsilon \quad (15)$$

and the termination criteria for the two remaining methods are similar. Alternatively, we also terminate if we exceed a time limit, as discussed in our numerical results. In the latter case, we evaluate the true cost of \mathbf{z}^t by computing its cost across each scenario in \mathcal{R} .

We remark that for some adversarial instances of Problem (1), using the same sample size at each iteration in conjunction with this termination criterion could lead to unattractive results where we terminate at a highly suboptimal solution with high probability (c.f. Morton 1998, Example 1). To address this issue and provide a confidence bound on our overall solution (accounting for multiple testing problems), we can increase the sample size at each iteration of the method in accordance with Morton (1998, Theorem 2) or use another sampling rule discussed therein (see also Bayraksan

and Morton 2011). However, owing to the single-tree implementation of our cutting-plane methods, as discussed in the next section, we do not test every candidate solution we generate when deciding to terminate. Therefore, as we observe in our numerical results, using the same sample size at each iteration is usually adequate. This is particularly true for the single-cut and k -cut methods, which, as discussed previously, often generate conservative lower bounds in practice, meaning that we often terminate at a computational time limit.

Finally, we remark that in circumstances where the total number of scenarios is relatively small, we can evaluate the true upper bound directly, rather than a stochastic estimate of the bound. Accordingly, we take this approach whenever the number of scenarios is sufficiently small.

3.2. Integrating Optimality Cuts Within a Branch-and-Cut Framework

Once our cut-generation and termination criterion schemes have been designed, they need to be embedded within a branch-and-cut framework in order to solve Problem (1) to certifiable optimality. Indeed, in the naive implementation of our algorithms described in the pseudocode of the previous section, we need to solve a mixed-integer problem at each iteration. For further scalability benefits, we can integrate our stochastic cut generation procedure within a state-of-the-art commercial mixed-integer solver (namely, **Gurobi** version 9.1.2) using **lazy constraint callbacks**, which accelerate cutting-plane methods by constructing a single branch-and-bound tree. For example, they have been used to implement deterministic cutting-plane algorithms in a highly efficient and relatively standard way; see, e.g., Fischetti et al. (2017, Section 4) .

Mixed-integer solvers assume that **lazy constraints** are binding at the point they are generated. Accordingly, they do not visit and do not generate **lazy constraints** twice at the same solution. Our stochastic cuts, however, are not binding, they provide a valid yet not necessarily tight lower bound. Therefore, when we implement our method with **lazy constraints**, the MIO solver can terminate with a highly suboptimal solution it deems optimal, because it (mistakenly) assumes the value of the cut generated at \mathbf{z}^t and evaluated at that point is precisely the cost of \mathbf{z}^t . To avoid this issue, we take a hybrid approach between single- and multi-tree branch-and-cut, which, to our knowledge, has not yet been described in the literature.

Namely, we maintain an outer loop where, at each iteration, we run a single-tree implementation of branch-and-cut with stochastic cutting-planes. During the branch-and-cut algorithm, we save all the cuts generated and imposed as **lazy constraints** within a separate cut pool. After the branch-and-cut algorithm, we randomly sample a subset of scenarios \mathcal{W} and compute the termination criterion described in the previous section to determine whether the solution returned by the branch-and-cut algorithm is indeed ϵ -optimal with high probability ($\alpha = 0.90$). By computing this convergence criterion at each iteration of the outer loop only, we mitigate the issue of multiple

hypothesis testing that would arise when testing the quality of a solution at each iteration of Algorithm 1–2 (inner loop). If the convergence criterion is met, we terminate the algorithm. Otherwise, we rerun the branch-and-cut algorithm and ensure the MIO solver no longer considers the previously generated `lazy constraints` as binding: We apply the constraints generated in the lazy cut pool as regular linear constraints, purge the lazy cut pool, and rerun the branch-and-cut algorithm. In addition to an optimality gap criterion, we also terminate the algorithm when the total computational time exceeds a predefined `TimeLimit`. We summarize this procedure in Algorithm 3. We remark that this approach is related to the notion of restarting a single-tree decomposition in a classical deterministic Benders scheme (see, e.g., Fischetti et al. 2016, Section 4.4).

Algorithm 3 Outer Loop for Stochastic Branch-and-Cut

- 1: **initialize** $CutPool = \emptyset$, $t = 0$
 - 2: **repeat**
 - 3: Increment $t \leftarrow t + 1$
 - 4: Initialize Algorithm 1/2 with constraints in $CutPool$.
 - 5: Run lazy-constraint implementation of Algorithm 1/2
 - 6: Save all lazy constraints generated in $CutPool$.
 - 7: Obtain candidate optimal solution \mathbf{z}^t .
 - 8: Obtain valid lower bound from the MIO branch-and-cut solver.
 - 9: Sample \mathcal{W} and compute $\bar{c}_{\alpha,t}$
 - 10: **until** (15) or `TimeLimit`
 - 11: Return \mathbf{z}^t , stochastic upper bound, and deterministic lower bound
-

Finally, in addition to the hybrid scheme described in this section, one could also consider a pure multi-tree implementation of our stochastic cutting-plane methods, as suggested in Section 2 and the classical network design literature (Geoffrion and Graves 1974). However, in preliminary numerical experiments, we found that such an approach is significantly slower because it involves solving a different MIO to generate each cut. Accordingly, we do not consider such an approach as part of our numerical experiments.

3.3. Accelerating the Convergence of our Approach

We now describe some practical enhancements to our stochastic cutting-plane approaches that improve their convergence, sometimes substantially; see also Fischetti et al. (2016, 2017), Bertsimas et al. (2021) for related discussions on accelerating the convergence of decomposition schemes.

Warm-Starting the Lower Bound: Cuts at the Root Node First, we warm-start our lower bound by applying cutting planes at the root node, as advocated by Fischetti et al. (2017), Bertsimas et al. (2021) among others. Specifically, we solve a Boolean relaxation of (3) using a continuous analog of our discrete cutting-plane method, and apply these cuts to the root node of our integer problem before running our branch-and-cut method. Note that this continuous cutting plane algorithm can also be implemented in a multi- or single-cut fashion, and in a deterministic or stochastic version. To balance the tightness of the formulation at the root node against the overall cost of computation, we impose a hard constraint on the total number of root node cuts applied (typically 10 or 20).

Warm-Starting the Upper Bound: We supply the initial network (without any new construction) as a warm-start for all methods we benchmark in our numerical experiments. This is guaranteed to be a feasible solution throughout our numerical experiments, due to our instance generation procedure. However, we do not implement a more sophisticated warm-starting strategy for any of our methods, in order to better understand the numerical behavior of our decomposition schemes.

We remark that, in practice, the Boolean relaxation could be randomly rounded to generate provably high-quality feasible solutions (c.f. Bertsimas et al. 2021, Section 3.2), and other heuristics specific to network design could be applied as well, as reviewed in the introduction.

4. Numerical Experiments

In this section, we numerically benchmark our sample-based Benders decomposition schemes on data-driven MCFND problems. We also compare their performance with their deterministic counterparts, and the performance of **Gurobi** on a perspective reformulation of the original MIO formulation (1). We consider a range of synthetic instances of varying sizes, including the so-called **R** instances commonly benchmarked in the literature (e.g., Crainic et al. 2016).

All experiments were conducted on MIT’s Supercloud Cluster (Reuther et al. 2018), which hosts Intel Xeon Platinum 8260 processors. All algorithms were implemented in Julia v1.7.3 (Bezanson et al. 2017) using JuMP v0.21.10 (Dunning et al. 2017) and Gurobi v9.1.2 (Gurobi Optimization, LLC 2022). The RAM memory allocated varies from 16GB to 140GB for the largest instances.

4.1. Problem Generation

We generate instances according to a methodology inspired by that of Günlük and Linderoth (2009) and Bertsimas et al. (2021). We construct a random graph by uniformly positioning $|\mathcal{N}|$ nodes over the unit square $[0, 1]^2$ and randomly sampling edges to construct a set of feasible edges \mathcal{E}_0 : We iteratively sample edges from the k -nearest neighbors graph (with $k = 6$) until we obtain a connected graph to ensure the feasibility of our instances. The construction cost for each edge, $c_{i,j}$ is drawn uniformly from $\mathcal{U}(1, 4)$. Each commodity $k \in \mathcal{K}$ corresponds to an all-to-one shortest

path problem with a single destination node $i_k \in \mathcal{N}$. For commodity k , we independently sample demands from all nodes $i' \in \mathcal{N}$, $d_{i'}^k$, uniformly between 5 to 20. We set $d_{i_k}^k := -\sum_{i' \neq i_k} d_{i'}^k$. We generate \mathcal{R} demand scenarios for each commodity accordingly. This process is repeated for every scenario $r \in \mathcal{R}$. Flow circulation costs, $f_{i,j}^k$, are proportional to the edge length (by a factor 10). The capacity of each arc is scaled based on the maximum cumulative demand across all scenarios: $B_{i,j} := \sum_{k \in \mathcal{K}} \sum_{(i,j) \in \mathcal{E}} \max_{r \in \mathcal{R}} d_{ij}^{k,r}$. Formally, we sample the capacity for arc (i,j) according to $u_{ij} \sim \mathcal{U}(1,4) \cdot B_{i,j}/|\mathcal{E}_0|$. We fix the cardinality constraint to $c_0 = 2|\mathcal{E}_0|$.

All in all, we generate instances with varying numbers of nodes $|\mathcal{N}|$, commodities $|\mathcal{K}|$, and scenarios $|\mathcal{R}|$, as described in Table 1. We later refer to these instances as small-, medium-, and large-scale instances based on the number of nodes $|\mathcal{N}|$.

Table 1 Dimensions of the MCFND problems generated, by scale (small-, medium-, and large-scale).

Scale	$ \mathcal{N} $	$ \mathcal{K} $	$ \mathcal{R} $
Small	{10,30,50,70}	× {5,10,25,50}	× {10,30,50,70,100}
Medium	{100,150,200}	× {5,10,25,50}	× {10,30,50,70,100}
Large	{300,500,700}	× {5,10,25,50}	× {10,30,50,70,100}

For our algorithms, we use two termination criteria: a time limit (3,600 or 7,200 seconds) and an optimality gap target $\epsilon = 1\%$ (with $\alpha = 0.90$ for our stochastic algorithms). Note that the time limit applies to the full outer-loop presented in Algorithm 3 (and not on each run of the branch-and-cut algorithm only). We also fix the regularization parameter γ to 1—we discuss its impact on our algorithms in Appendix B. We warm-start all methods with the original connected graph as an initial solution.

4.2. Comparison of Different Stochastic Cutting-Plane Algorithms

In this section, we benchmark the variants of the stochastic cutting plane algorithm proposed in Section 2, namely the multi-, single-, and k -cut ($k = |\mathcal{K}|/2$) algorithms, in terms of their ability to obtain a certifiably near-optimal solution with high confidence. We also measure the impact of warm-starting these methods with cuts obtained from solving the perspective relaxation with a multi- or single-cut stochastic cutting plane algorithm, and applying these cuts at the root node in our branch-and-cut scheme (which we refer to as multi-cut or single-cut root node cuts respectively). For all stochastic methods, we use a sampling rate, $|\mathcal{R}_t|/|\mathcal{R}|$, of 10% and a time limit of 3,600 seconds. Accordingly, we report average computational time (capped at 3,600 seconds) for solving our small and medium-scale synthetic instances in Table 2. To augment these results, Table 3 reports the average optimality gap at termination, and Table C.1 (see Appendix C.1) reports the

Table 2 Computational time (in seconds) of the multi-, single-, and k -cut stochastic cutting plane algorithm, with different warm-start strategies at the root node (none, multi-cut, and single-cut root node cuts). Metrics are averaged across instances with the same number of nodes $|\mathcal{N}|$.

$ \mathcal{N} $	Multi-Cut			Single-Cut			k -Cut		
	None	Multi	Single	None	Multi	Single	None	Multi	Single
10	175.7	113.1	59.5	85.9	137.6	69.0	124.0	93.1	82.5
30	2587.8	3013.8	2905.4	2773.9	2357.9	1741.8	2582.8	2131.4	2227.7
50	3073.6	3600.0	3600.0	2542.7	2477.4	1715.9	2727.3	2098.1	2532.2
70	3170.6	3527.2	3485.9	2786.9	2956.4	2729.8	2753.3	2590.7	2612.1
100	3220.4	3600.0	3600.0	2849.6	2881.5	2732.6	2992.2	2625.0	3138.7
150	3600.0	3600.0	3547.9	3219.1	2883.1	2779.8	3091.4	2758.7	3113.1
200	3600.0	3600.0	3600.0	2884.2	3093.1	2702.7	3231.0	2761.6	2961.4

Table 3 Relative optimality gap (in %) at termination for multi-, single-, and k -cut algorithms, with different warm-start strategies at the root node. Metrics are averaged across instances with same number of nodes $|\mathcal{N}|$.

$ \mathcal{N} $	Multi-Cut			Single-Cut			k -Cut		
	None	Multi	Single	None	Multi	Single	None	Multi	Single
10	0.20	0.05	0.93	0.13	0.06	0.23	1.52	0.03	0.27
30	9.66	5.26	4.21	12.03	10.84	4.29	10.71	4.99	4.56
50	11.08	2.78	2.52	7.98	8.60	2.63	7.33	2.74	3.01
70	35.45	10.71	5.10	20.44	17.17	6.13	31.03	11.98	9.17
100	37.05	10.40	4.25	22.69	17.58	6.41	33.99	15.62	8.04
150	46.32	12.18	4.53	30.39	20.72	9.59	38.82	13.46	10.19
200	46.27	17.43	14.18	33.94	22.44	10.36	43.54	15.50	11.75

fraction of instances solved within the time limit. Note that the optimality gaps reported in Table 3 are computed using the true cost of the incumbent solution, using all scenarios in \mathcal{R} .

We observe that both the multi-cut and single-cut warm-start strategies are very effective in reducing the relative optimality gap at termination. Indeed, our root node strategies halve the optimality gap at termination compared to not applying cuts at the root node. However, a single-cut strategy at the root node appears to outperform a multi-cut root node strategy in terms of the relative gap at termination. All in all, applying a single-cut strategy warm-started with a single-cut method at the root node appears to be the most numerically efficient strategy, because this strategy usually terminates faster than any other strategy and has a relative optimality gap comparable to or better than any other strategy. Yet, multi-cut with a single-cut root node strategy and k -cut with a single-cut root node strategy are also competitive, and indeed multi-cut obtains a smaller average optimality gap at termination at the price of a higher average runtime.

Next, we investigate the number of iterations of the outer loop performed by our methods; recall that in Section 3, we proposed an outer loop procedure that allows our sampling approach to be safely integrated within a branch-and-cut procedure, without requiring a new branch-and-bound tree each time we generate a cut. To this end, Figure 1 depicts the number of outer-loop iterations performed by our single-cut algorithm on the small- and medium-scale instances. We observe that only one iteration of the outer loop is performed in many cases ($\sim 50\%$ for small-scale and 70% for medium-scale instances). However, in the remaining cases, the first iteration of branch-and-cut with stochastic cuts terminates with a solution that is not ϵ -optimal but Algorithm 3 is very efficient, requiring a limited number of additional iterations to identify an optimal solution. This verifies that a single outer loop iteration often wrongly terminates at a solution that is not optimal. On the other hand, only a small number of iterations of the outer loop are usually needed to achieve optimality. Therefore, the tractability of our approach is not compromised by the outer loop.

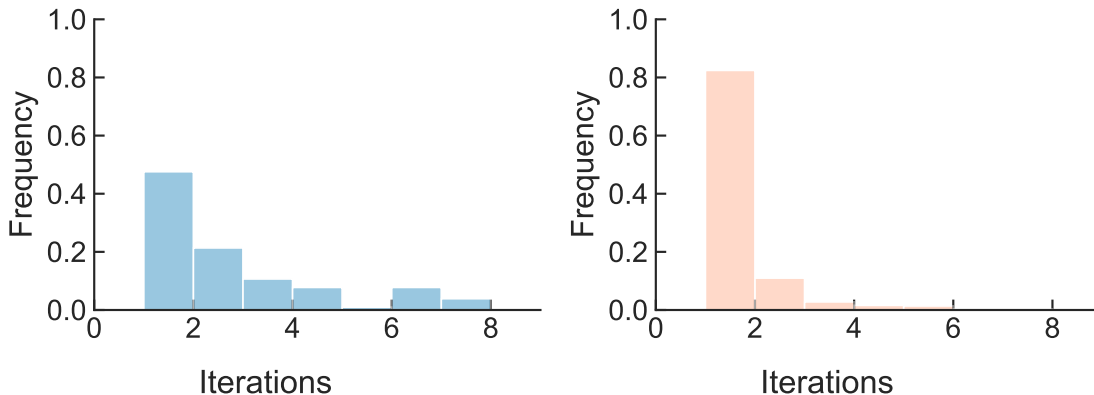


Figure 1 Distribution of the number of outer-loop iterations required by the single- and k -cut stochastic cutting plane algorithms with single-cut root node cuts on small-scale (left) and medium-scale (right) instances.

4.3. Benchmarking Scalability on Synthetic Instances

We now compare the performance of our stochastic cutting plane methods (single- and k -cut) against two benchmarks: (a) solving Problem (1)’s perspective reformulation directly with `Gurobi`, and (b) a deterministic single-cut method. For our stochastic approaches, we use a sampling rate of 10%. We impose a time limit of 7,200 seconds for all methods. To calibrate our approaches and verify their correctness, we use the smallest instances to verify that all methods terminate with the same optimal solution (see Table C.4 in Appendix C.2).

We report the average computational time and optimality gap of all methods, on the small-, medium-, and large-scale instances, in Table 4, with metrics averaged over instances with the same number of nodes $|\mathcal{N}|$.

Table 4 Runtime (in seconds) and final optimality gap (in %) for each algorithm, averaged over instances with the same number of nodes $|\mathcal{N}|$.

$ \mathcal{N} $	Gurobi with (1)		Deterministic		Stochastic		Stochastic k -Cut	
	Runtime	Gap	Runtime	Gap	Runtime	Gap	Runtime	Gap
10	363.0	0.00	303.6	10.06	83.0	0.27	69.7	1.10
30	7200.0	38.06	7200.0	9.80	4496.1	5.29	4583.7	4.50
50	7200.0	63.42	7192.4	7.70	4166.6	3.02	5813.2	2.79
70	7200.0	74.29	7200.0	10.17	4951.0	5.64	5007.6	5.39
100	-	-	6899.5	13.02	5168.7	4.94	6196.7	4.62
150	-	-	7004.4	12.54	6140.7	4.33	5822.1	5.31
200	-	-	6389.0	16.42	5715.2	5.32	5674.7	4.25
300	-	-	-	-	5025.8	9.48	4866.5	8.80
500	-	-	-	-	5521.4	19.01	6380.5	26.82
700	-	-	-	-	5420.5	37.32	5734.6	40.72

We observe that a perspective reformulation of the original formulation (1) cannot be solved by Gurobi with 100 or more nodes within the time (2 hours) and memory ($> 72\text{GB}$) limits. Indeed, while this approach converges within minutes for instances with ten nodes, it fails to identify an optimal solution within the two-hour time limit for instances with 20-70 nodes and terminates with large optimality gaps ($> 30\%$) on average. On the other hand, a deterministic Benders decomposition scheme reaches optimality gaps that are an order of magnitude smaller on instances with 20-70 nodes, scales to instances with up to 200 nodes, but fails to recover a solution with a meaningful optimality gap within the time limit for larger problems.

Our stochastic cutting plane algorithms significantly improve upon their deterministic counterpart. On small- and medium-scale instances, they reduce the average computational time by 30-40% on the small instances and 10-25% on the medium ones. A comparison in terms of average computational times might be misleading, however, because of the time limit, and because many of these instances are not solved to ϵ -optimality. Accordingly, we also compare in terms of the optimality gap. We observe that our stochastic cutting-plane algorithms terminate with gaps half the size of deterministic algorithms (i.e., around 5% for the instances with 70-200 nodes compared with 10-15% for the deterministic approach). Finally, the stochastic cutting-plane methods successfully scale to instances with 300-700 nodes. They obtain bound gaps of less than 10% for instances with 300 nodes and less than 40% with 700 nodes (within the same 2-hour time limit), which is one order of magnitude larger than the instances solved in the stochastic network design literature (Crainic et al. 2021a). Figure C.1 in Appendix C.2 displays the optimality gap achieved for each values of $|\mathcal{N}|$, $|\mathcal{R}|$, and $|\mathcal{K}|$, and shows our method is most sensitive to the number of commodities and nodes.

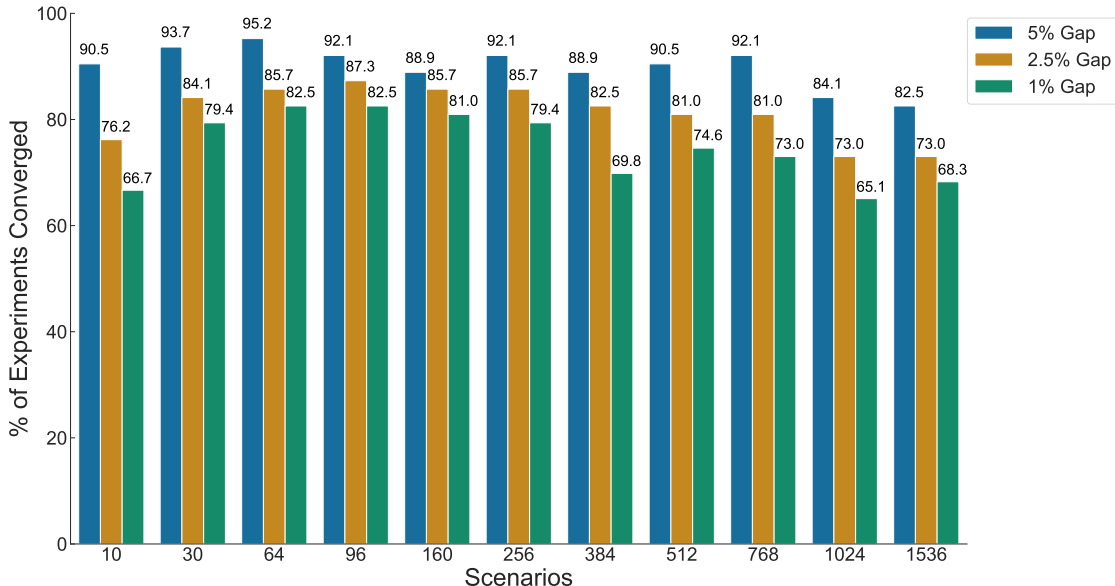


Figure 2 Fraction of instances from Crainic et al. (2000) solved by the single-cut stochastic cutting plane algorithm within a 2-hour time limit and for different target optimality gaps ϵ .

4.4. Benchmarking on the Instances from Crainic et al. (2000)

We now benchmark our methods on network design instances from the literature. Namely, we use the so-called **R** instances, originally introduced by Crainic et al. (2000) for deterministic network design problems and widely used to assess the scalability of methods with respect to the number of scenarios (Crainic et al. 2021b, 2016, Boland et al. 2016). We use a total 63 instances (see details in Appendix C.3) with 10–20 nodes and 10–50 commodities. So, compared to our synthetic instances, the ratio $|\mathcal{K}|/|\mathcal{N}|$ is higher for these instances. We generate instances with the same structure (nodes, arcs, commodities, capacities, and costs) as the nominal **R** instances, and with demand scenarios that are random perturbations of the nominal demand. The number of scenarios generated varies from 10 to 1,536. We evaluate the performance of our single-cut stochastic cutting plane algorithm with a 2% sampling rate and a time limit of 7,200 seconds for all methods.

Figure 2 represents the proportion of instances we solve to a different optimality gaps ϵ (we keep $\alpha = 0.90$), as the number of scenarios $|\mathcal{R}|$ increases. Our results indicate that the fraction of instances solved to optimality within 2 hours seems to generally decrease as the number of scenarios increases. However, we observe that we can achieve an optimality gap below 2.5% for the majority of instances ($\geq 73\%$) across all values of $|\mathcal{R}|$.

Table 5 reports the average optimality gap achieved after 7,200 seconds, as the number of scenarios increases. Although it increases slightly with the number of scenarios, we observe that the average optimality gap remains around 2.5%. In comparison, Crainic et al. (2021b) only provides valid lower bounds for up to 256 scenarios. We acknowledge, however, that our instance generation

strategy differs slightly from theirs, and our formulation involves a strongly convex regularization term in the objective, so an apples-to-apples comparison between our results is not possible.

Table 5 Optimality gap (in %) achieved by the single-cut stochastic cutting-plane algorithm with a 2-hour time limit, as the same number of scenarios $|\mathcal{R}|$ generated for the instances in Crainic et al. (2000) increases.

Number of scenarios	10	30	64	96	160	256	384	512	768	1,024	1,536
Gap	1.73	1.6	1.6	1.42	1.52	1.46	1.72	1.84	1.58	1.98	2.21

Finally, Figure 3 represents the average runtime (left panel) and the number of cuts (in log terms, right panel) of our stochastic cutting-plane algorithm for varying sampling rates $\mathcal{R}_t/|\mathcal{R}|$ and different number of scenarios $|\mathcal{R}|$. A single-cut root node cut routine, which samples at a 25% rate, is used to initiate all algorithms. As expected, we observe that the number of cuts required for our algorithm to converge increases as the number of scenarios sampled for each cut decreases. However, this loss in efficiency in number of cuts is outweighed by the computational benefit when generating each cut, so a lower sampling rate leads to an overall lower total computational time.

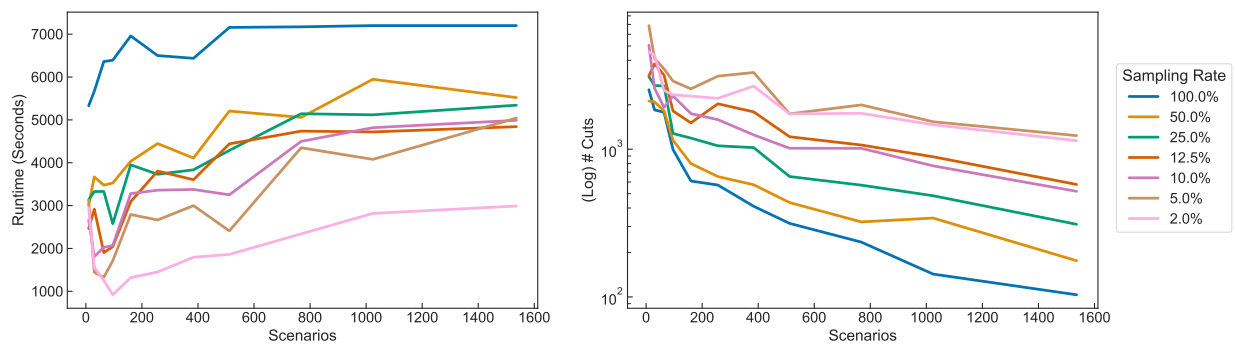


Figure 3 Computational time (left panel) and log-number of cuts generated (right panel) of the single-cut stochastic cutting plane algorithm on the R instances, for different sampling rates.

5. Conclusion

We propose a stochastic Benders decomposition scheme which solves large-scale data-driven stochastic network design problems. Our approach mitigates the high computational cost of generating each cut by sampling a subset of the data at each iteration, while applying a dual-averaging technique to ensure that the cuts generated remain valid for the original problem. We also propose an outer loop technique to ensure the safe termination of our algorithm when the Benders decomposition scheme is implemented via lazy callbacks. We consider multi-, single-, and k -cut variants of our algorithm and discuss its implementation within a branch-and-cut solver. In numerical experiments, we demonstrate that our stochastic decomposition schemes obtain optimality gaps of 4–5%

on instances with 100–200 nodes, compared to 13–16% for deterministic Benders schemes. Moreover, we obtain bound gaps of 40% on instances with up to 700 nodes and 50 commodities, i.e., problem sizes an order of magnitude larger than any instances addressed by exact methods in the literature. Beyond network design, we believe our approach could be applied to other two-stage stochastic optimization problems addressed via sample average approximations.

References

- Agrawal A, Klein P, Ravi R (1991) When trees collide: An approximation algorithm for the generalized Steiner problem on networks. *Proceedings of ACM Symposium on Theory of Computing*, 134–144.
- Atamtürk A, Günlük O (2018) A note on capacity models for network design. *Operations Research Letters* 46(4):414–417.
- Atamtürk A, Günlük O (2021) Multicommodity multifacility network design. *Network Design with Applications to Transportation and Logistics*, 141–166 (Springer).
- Balakrishnan A, Magnanti TL, Shulman A, Wong RT (1991) Models for planning capacity expansion in local access telecommunication networks. *Annals of Operations Research* 33(4):237–284.
- Bardenet R, Maillard OA (2015) Concentration inequalities for sampling without replacement. *Bernoulli* 21(3):1361–1385.
- Barnhart C, Belobaba P, Odoni AR (2003) Applications of operations research in the air transport industry. *Transportation Science* 37(4):368–391.
- Bayraksan G, Morton DP (2011) A sequential sampling procedure for stochastic programming. *Operations Research* 59(4):898–913.
- Beale EM (1955) On minimizing a convex function subject to linear inequalities. *Journal of the Royal Statistical Society: Series B (Methodological)* 17(2):173–184.
- Bertsekas DP (1999) *Nonlinear Optimization* (Athena Scientific, Belmont).
- Bertsimas D, Cory-Wright R, Pauphilet J (2021) A unified approach to mixed-integer optimization problems with logical constraints. *SIAM Journal on Optimization* 31(3):2340–2367.
- Bertsimas D, Li ML (2022) Stochastic cutting planes for data-driven optimization. *INFORMS Journal on Computing* 34(5):2400–2409.
- Bertsimas D, Teo CP (1998) From valid inequalities to heuristics: A unified view of primal-dual approximation algorithms in covering problems. *Operations Research* 46(4):503–514.
- Bertsimas D, Tsitsiklis JN (1997) *Introduction to Linear Optimization*, volume 6 (Athena Scientific Belmont, MA).
- Bezanson J, Edelman A, Karpinski S, Shah VB (2017) Julia: A fresh approach to numerical computing. *SIAM Review* 59(1):65–98.

- Bienstock D, Chopra S, Günlük O, Tsai CY (1998) Minimum cost capacity installation for multicommodity network flows. *Mathematical Programming* 81(2):177–199.
- Binato S, Pereira MVF, Granville S (2001) A new Benders decomposition approach to solve power transmission network design problems. *IEEE Transactions on Power Systems* 16(2):235–240.
- Birge JR, Louveaux F (2011) *Introduction to Stochastic Programming*. Springer Series in Operations Research and Financial Engineering (New York, NY: Springer New York).
- Birge JR, Louveaux FV (1988) A multicut algorithm for two-stage stochastic linear programs. *European Journal of Operational Research* 34(3):384–392.
- Bixby RE (2012) A brief history of linear and mixed-integer programming computation. *Documenta Mathematica* 2012:107–121.
- Boland N, Fischetti M, Monaci M, Savelsbergh M (2016) Proximity Benders: a decomposition heuristic for stochastic programs. *Journal of Heuristics* 22(2):181–198.
- Boyce DE, Farhi A, Weischedel R (1973) Optimal network problem: a branch-and-bound algorithm. *Environment and Planning A* 5(4):519–533.
- Ceria S, Soares J (1999) Convex programming for disjunctive convex optimization. *Mathematical Programming* 86(3):595–614.
- Contreras I, Cordeau JF, Laporte G (2011) Benders decomposition for large-scale uncapacitated hub location. *Operations Research* 59(6):1477–1490.
- Cornuejols G, Nemhauser GL, Wolsey LA (1980) A canonical representation of simple plant location problems and its applications. *SIAM Journal on Algebraic Discrete Methods* 1(3):261–272.
- Costa AM (2005) A survey on Benders decomposition applied to fixed-charge network design problems. *Computers & Operations Research* 32(6):1429–1450.
- Crainic TG, Gendreau M, Farvolden JM (2000) A simplex-based Tabu search method for capacitated network design. *INFORMS Journal on Computing* 12(3):223–236.
- Crainic TG, Gendreau M, Gendron B, eds. (2021a) *Network Design with Applications to Transportation and Logistics* (Cham: Springer International Publishing).
- Crainic TG, Hewitt M, Maggioni F, Rei W (2021b) Partial Benders decomposition: general methodology and application to stochastic network design. *Transportation Science* 55(2):414–435.
- Crainic TG, Rei W, Hewitt M, Maggioni F (2016) *Partial Benders Decomposition Strategies for Two-Stage Stochastic Integer Programs*, volume 37 (CIRRELT).
- Dantzig GB (1955) Linear programming under uncertainty. *Management Science* 1(3-4):197–206.
- Dantzig GB, Infanger G (1993) Multi-stage stochastic linear programs for portfolio optimization. *Annals of Operations Research* 45(1):59–76.

-
- Davis D, Drusvyatskiy D, Kakade S, Lee JD (2020) Stochastic subgradient method converges on tame functions. *Foundations of Computational Mathematics* 20(1):119–154.
- de Camargo RS, Miranda Jr G, Luna HP (2008) Benders decomposition for the uncapacitated multiple allocation hub location problem. *Computers & operations research* 35(4):1047–1064.
- De Matos VL, Philpott AB, Finardi EC (2015) Improving the performance of stochastic dual dynamic programming. *Journal of Computational and Applied Mathematics* 290:196–208.
- Dunning I, Huchette J, Lubin M (2017) JuMP: A modeling language for mathematical optimization. *SIAM Review* 59(2):295–320.
- Erickson RE, Monma CL, Veinott Jr AF (1987) Send-and-split method for minimum-concave-cost network flows. *Mathematics of Operations Research* 12(4):634–664.
- Fábián CI (2000) Bundle-type methods for inexact data. *Central European Journal of Operations Research* 8(1):35–55.
- Fischetti M, Ljubić I, Sinnl M (2016) Benders decomposition without separability: A computational study for capacitated facility location problems. *European Journal of Operational Research* 253(3):557–569.
- Fischetti M, Ljubić I, Sinnl M (2017) Redesigning Benders Decomposition for Large-Scale Facility Location. *Management Science* 63(7):2146–2162.
- Florian M, Bushell G, Ferland J, Guerin G, Nastansky L (1976) The engine scheduling problem in a railway network. *INFOR: Information Systems and Operational Research* 14(2):121–138.
- Frangioni A, Gentile C (2006) Perspective cuts for a class of convex 0–1 mixed integer programs. *Mathematical Programming* 106(2):225–236.
- Frangioni A, Gentile C (2009) A computational comparison of reformulations of the perspective relaxation: SOCP vs. cutting planes. *Operations Research Letters* 37(3):206–210.
- Garey MR, Johnson DS (1977) The rectilinear steiner tree problem is np-complete. *SIAM Journal on Applied Mathematics* 32(4):826–834.
- Gendron B, Crainic TG, Frangioni A (1999) Multicommodity capacitated network design. *Telecommunications Network Planning*, 1–19 (Springer).
- Gendron B, Hanafi S, Todosijević R (2018) Matheuristics based on iterative linear programming and slope scaling for multicommodity capacitated fixed charge network design. *European Journal of Operational Research* 268(1):70–81.
- Geoffrion AM (1972) Generalized Benders decomposition. *Journal of Optimization Theory and Applications* 10(4):237–260.
- Geoffrion AM, Graves GW (1974) Multicommodity distribution system design by Benders decomposition. *Management Science* 20(5):822–844.

-
- Goemans MX, Bertsimas DJ (1993) Survivable networks, linear programming relaxations and the parsimonious property. *Mathematical Programming* 60(1):145–166.
- Grimmett G, Stirzaker D (2020) *Probability and Random Processes* (Oxford university press).
- Guigues V (2020) Inexact cuts in stochastic dual dynamic programming. *SIAM Journal on Optimization* 30(1):407–438.
- Günlük O (1999) A branch-and-cut algorithm for capacitated network design problems. *Mathematical Programming* 86(1):17–39.
- Günlük O, Linderoth J (2009) Perspective reformulations of mixed integer nonlinear programs with indicator variables. *Mathematical Programming* 183–205.
- Gurobi Optimization, LLC (2022) Gurobi Optimizer Reference Manual.
- Higle JL, Sen S (1991) Stochastic decomposition: An algorithm for two-stage linear programs with recourse. *Mathematics of Operations Research* 16(3):650–669.
- Higle JL, Sen S (1996a) Duality and statistical tests of optimality for two stage stochastic programs. *Mathematical Programming* 75(2):257–275.
- Higle JL, Sen S (1996b) *Stochastic Decomposition: A Statistical Method for Large Scale Stochastic Linear Programming*, volume 8 (Springer Science & Business Media).
- Infanger G (1992) Monte Carlo (importance) sampling within a Benders decomposition algorithm for stochastic linear programs. *Annals of Operations Research* 39(1):69–95.
- Ke Q, Ferrara E, Radicchi F, Flammini A (2015) Defining and identifying sleeping beauties in science. *Proceedings of the National Academy of Sciences* 112(24):7426–7431.
- Kingma DP, Ba J (2014) Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* .
- Magnanti TL, Mirchandani P, Vachani R (1993) The convex hull of two core capacitated network design problems. *Mathematical Programming* 60(1):233–250.
- Magnanti TL, Mirchandani P, Vachani R (1995) Modeling and solving the two-facility capacitated network loading problem. *Operations Research* 43(1):142–157.
- Magnanti TL, Wong RT (1981) Accelerating Benders decomposition: Algorithmic enhancement and model selection criteria. *Operations Research* 29(3):464–484.
- Magnanti TL, Wong RT (1984) Network design and transportation planning: models and algorithms. *Transportation Science* 18(1):1–55.
- Mak WK, Morton DP, Wood RK (1999) Monte carlo bounding techniques for determining solution quality in stochastic programs. *Operations Research Letters* 24(1-2):47–56.
- Morton DP (1998) Stopping rules for a class of sampling-based stochastic programming algorithms. *Operations Research* 46(5):710–718.

- Pereira MV, Pinto LM (1991) Multi-stage stochastic optimization applied to energy planning. *Mathematical Programming* 52(1):359–375.
- Pishvae MS, Razmi J, Torabi SA (2014) An accelerated Benders decomposition algorithm for sustainable supply chain network design under uncertainty: A case study of medical needle and syringe supply chain. *Transportation Research Part E: Logistics and Transportation Review* 67:14–38.
- Rahmaniani R, Crainic TG, Gendreau M, Rei W (2018) Accelerating the Benders decomposition method: Application to stochastic network design problems. *SIAM Journal on Optimization* 28(1):875–903.
- Rei W, Cordeau JF, Gendreau M, Soriano P (2009) Accelerating Benders decomposition by local branching. *INFORMS Journal on Computing* 21(2):333–345.
- Reuther A, Kepner J, Byun C, Samsi S, Arcand W, Bestor D, Bergeron B, Gadepally V, Houle M, Hubbell M, Jones M, Klein A, Milechin L, Mullen J, Prout A, Rosa A, Yee C, Michaleas P (2018) Interactive supercomputing on 40,000 cores for machine learning and data analysis. *2018 IEEE High Performance extreme Computing Conference (HPEC)*, 1–6 (IEEE).
- Richardson R (1976) An optimization approach to routing aircraft. *Transportation Science* 10(1):52–71.
- Rodríguez-Martín I, Salazar-González JJ (2010) A local branching heuristic for the capacitated fixed-charge network design problem. *Computers & Operations Research* 37(3):575–581.
- Santoso T, Ahmed S, Goetschalckx M, Shapiro A (2005) A stochastic programming approach for supply chain network design under uncertainty. *European Journal of Operational Research* 167(1):96–115.
- Schmidt M, Le Roux N, Bach F (2017) Minimizing finite sums with the stochastic average gradient. *Mathematical Programming* 162(1):83–112.
- Shapiro A, Dentcheva D, Ruszczyński A (2021) *Lectures on Stochastic Programming: Modeling and Theory* (SIAM).
- Smith JE, Winkler RL (2006) The optimizer’s curse: Skepticism and postdecision surprise in decision analysis. *Management Science* 52(3):311–322.
- Stubbs RA, Mehrotra S (1999) A branch-and-cut method for 0-1 mixed convex programming. *Mathematical Programming* 86(3):515–532.
- Trukhanov S, Ntaimo L, Schaefer A (2010) Adaptive multicut aggregation for two-stage stochastic linear programs with recourse. *European Journal of Operational Research* 206(2):395–406.
- Van Roy TJ, Wolsey LA (1985) Valid inequalities and separation for uncapacitated fixed charge networks. *Operations Research Letters* 4(3):105–112.
- Van Slyke RM, Wets R (1969) L-shaped linear programs with applications to optimal control and stochastic programming. *SIAM Journal on Applied Mathematics* 17(4):638–663.
- Wets RJB (1966) Programming under uncertainty: the equivalent convex program. *SIAM Journal on Applied Mathematics* 14(1):89–105.

-
- Xie W, Deng X (2020) Scalable algorithms for the sparse ridge regression. *SIAM Journal on Optimization* 30(4):3359–3386.
- You F, Grossmann IE (2013) Multicut Benders decomposition algorithm for process supply chain planning under uncertainty. *Annals of Operations Research* 210(1):191–211.
- Zakeri G, Philpott AB, Ryan DM (2000) Inexact cuts in Benders decomposition. *SIAM Journal on Optimization* 10(3):643–657.

Appendix A: Omitted Proofs

A.1. Proof of Proposition 1

Proof of Proposition 1 The minimization problem defining $f(\mathbf{z}; \mathbf{d})$ can be seen as the sum of two minimization problems

$$\min_{\mathbf{x}^k \in \mathbb{R}_+^{\mathcal{E}}, k \in \mathcal{K}} \sum_{k \in \mathcal{K}} \langle \mathbf{f}^k, \mathbf{x}^k \rangle \text{ s.t. } \mathbf{A}\mathbf{x}^k = \mathbf{d}^k, \forall k \in \mathcal{K},$$

and

$$\min_{\mathbf{y} \in \mathbb{R}^{\mathcal{E}}} \frac{1}{2\gamma} \sum_{(i,j) \in \mathcal{E}} y_{i,j}^2 \text{ s.t. } y_{i,j} \leq u_{i,j}, \forall (i,j) \in \mathcal{E},$$

$$y_{i,j} = 0 \text{ if } z_{i,j} = 0, \forall (i,j) \in \mathcal{E},$$

coupled via the constraints $\sum_{k \in \mathcal{K}} x_{i,j}^k = y_{i,j}, \forall (i,j) \in \mathcal{E}$. Therefore, by associating a dual variable $\alpha_{i,j} \in \mathbb{R}$ with each coupling constraint, we rewrite $f(\mathbf{z}; \mathbf{d})$ as

$$\min_{\substack{\mathbf{x}^k \in \mathbb{R}_+^{\mathcal{E}}, k \in \mathcal{K}: \\ \mathbf{A}\mathbf{x}^k = \mathbf{d}^k, \forall k \in \mathcal{K}}} \min_{\substack{\mathbf{y} \in \mathbb{R}^{\mathcal{E}}: \\ y_{i,j} \leq u_{i,j}, \forall (i,j) \in \mathcal{E} \\ y_{i,j} = 0 \text{ if } z_{i,j} = 0, \forall (i,j) \in \mathcal{E}}} \max_{\alpha \in \mathbb{R}_+^{\mathcal{E}}} \sum_{k \in \mathcal{K}} \langle \mathbf{f}^k - \alpha, \mathbf{x}^k \rangle + \sum_{(i,j) \in \mathcal{E}} \left(\alpha_{i,j} y_{i,j} + \frac{1}{2\gamma} y_{i,j}^2 \right).$$

By invoking standard results on saddle-point theorems (see, e.g., Bertsekas 1999), the order of the minimization and maximization operators on the function $f(\mathbf{z}, \mathbf{d})$ can be exchanged² without altering the objective value. Moreover, after exchanging these operators, we can compute the dual of each minimization problem separately. Indeed,

$$\min_{\substack{\mathbf{x}^k \in \mathbb{R}_+^{\mathcal{E}}, k \in \mathcal{K}: \\ \mathbf{A}\mathbf{x}^k = \mathbf{d}^k, \forall k \in \mathcal{K}}} \sum_{k \in \mathcal{K}} \langle \mathbf{f}^k - \alpha, \mathbf{x}^k \rangle = \max_{\substack{\mathbf{p}^k \in \mathbb{R}^{\mathcal{N}}: \\ \mathbf{A}^\top \mathbf{p}^k \leq \mathbf{f}^k - \alpha, \forall k \in \mathcal{K}}} \sum_{k \in \mathcal{K}} \langle \mathbf{p}^k, \mathbf{d}^k \rangle.$$

Second, to dualize

$$\min_{\mathbf{y} \in \mathbb{R}^{\mathcal{E}}} \sum_{(i,j) \in \mathcal{E}} \left(\alpha_{i,j} y_{i,j} + \frac{1}{2\gamma} y_{i,j}^2 \right) \text{ s.t. } y_{i,j} \leq u_{i,j}, \forall (i,j) \in \mathcal{E},$$

$$y_{i,j} = 0 \text{ if } z_{i,j} = 0, \forall (i,j) \in \mathcal{E},$$

let us first observe that we can omit the logical constraints by considering the change of variables $y_{i,j} = z_{i,j} w_{i,j}$ for $\mathbf{w} \in \mathbb{R}^{\mathcal{E}}$. Hence, we obtain

$$\min_{\mathbf{w} \in \mathbb{R}^{\mathcal{E}}} \sum_{(i,j) \in \mathcal{E}} \left[z_{i,j} \alpha_{i,j} w_{i,j} + \frac{1}{2\gamma} z_{i,j} w_{i,j}^2 \right] \text{ s.t. } z_{i,j} w_{i,j} \leq u_{i,j}, \forall (i,j) \in \mathcal{E}$$

$$= \max_{\beta \in \mathbb{R}_+^{\mathcal{E}}} \min_{\mathbf{u} \in \mathbb{R}^{\mathcal{E}}} -\langle \beta, \mathbf{u} \rangle + \sum_{(i,j) \in \mathcal{E}} \left[z_{i,j} (\alpha_{i,j} + \beta_{i,j}) w_{i,j} + \frac{1}{2\gamma} z_{i,j} w_{i,j}^2 \right]$$

$$= \max_{\beta \in \mathbb{R}_+^{\mathcal{E}}} -\langle \beta, \mathbf{u} \rangle - \frac{\gamma}{2} \sum_{(i,j) \in \mathcal{E}} z_{i,j} (\alpha_{i,j} + \beta_{i,j})^2.$$

All together, we obtain the desired reformulation. \square

² In general, we require that a constraint qualification holds to be able to exchange the order of minimization and maximization operators (see, e.g., Bertsekas 1999). However, all constraints in Problem (2) are linear and it has a convex quadratic objective. Therefore, we can exchange the order of the operators in an assumption-free manner.

A.2. Proof of Proposition 2

In this section, we provide a proof of Proposition 2. To clarify the presentation, we adopt a lighter set of notations:

Fix \mathbf{z} . For any $r \in \mathcal{R}$, we denote $\boldsymbol{\xi}^r := (\boldsymbol{\alpha}^r, \boldsymbol{\beta}^r, \mathbf{p}^r)$ the optimal dual solutions of (4) for $\mathbf{d} = \mathbf{d}^r$. For any subset $\mathcal{S} \subseteq \mathcal{R}$, let us denote $\bar{\boldsymbol{\xi}}^{\mathcal{S}} := \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{S}} \boldsymbol{\xi}^r$ the average of the optimal dual solutions $\boldsymbol{\xi}^r$ for $r \in \mathcal{S}$. For a random $\mathcal{S} \subseteq \mathcal{R}$ of fixed size $|\mathcal{S}|$, we will analyze the sub-optimality gap of $\bar{\boldsymbol{\xi}}^{\mathcal{S}}$, i.e., the quantity $q(\mathbf{z}, \boldsymbol{\xi}^r; \mathbf{d}^r) - q(\mathbf{z}, \bar{\boldsymbol{\xi}}^{\mathcal{S}}; \mathbf{d}^r) (\geq 0)$, for scenarios $r \in \mathcal{S}^c := \mathcal{R} \setminus \mathcal{S}$.

Proof of Proposition 2 Let us denote $M := \max_{r \in \mathcal{R}} \|\boldsymbol{\xi}^r\|_{\infty}$. Since $\|\boldsymbol{\xi}^r\|_{\infty} \leq M$, then $\|\bar{\boldsymbol{\xi}}^{\mathcal{S}}\|_{\infty} \leq M$ and there exists some constant $L > 0$ such that, for any \mathcal{S} and any $r \notin \mathcal{S}$

$$\left| q(\mathbf{z}, \bar{\boldsymbol{\xi}}^{\mathcal{S}}, \mathbf{d}^r) - q(\mathbf{z}, \boldsymbol{\xi}^r, \mathbf{d}^r) \right| \leq L \|\bar{\boldsymbol{\xi}}^{\mathcal{S}} - \boldsymbol{\xi}^r\|.$$

We further decompose the right-hand side via a triangle inequality and sum the inequalities above across all $r \notin \mathcal{S}$ to obtain

$$\sum_{r \in \mathcal{S}^c} \left| q(\mathbf{z}, \bar{\boldsymbol{\xi}}^{\mathcal{S}}, \mathbf{d}^r) - q(\mathbf{z}, \boldsymbol{\xi}^r, \mathbf{d}^r) \right| \leq L |\mathcal{S}^c| \|\bar{\boldsymbol{\xi}}^{\mathcal{S}} - \bar{\boldsymbol{\xi}}^{\mathcal{R}}\| + L \sum_{r \in \mathcal{S}^c} \|\bar{\boldsymbol{\xi}}^{\mathcal{R}} - \boldsymbol{\xi}^r\|.$$

The first term corresponds to the difference between $\bar{\boldsymbol{\xi}}^{\mathcal{R}}$ and an unbiased estimate obtained via sampling without replacement. Denote d the dimension of $\boldsymbol{\xi}$. Hence, since the $\boldsymbol{\xi}^r$ are uniformly bounded, by Bardenet and Maillard (2015, corollary 2.5), there exists some universal constant M_1 such that for any $\delta > 0$, we have, with probability $1 - \delta$ on the subset \mathcal{S} of fixed size $|\mathcal{S}|$,

$$\|\bar{\boldsymbol{\xi}}^{\mathcal{S}} - \bar{\boldsymbol{\xi}}^{\mathcal{R}}\| \leq M_1 \sqrt{d \left(\frac{1}{|\mathcal{S}|} - \frac{1}{|\mathcal{R}|} \right) \log(1/\delta)}, \quad (16)$$

For the second term, we simply use the bound

$$\sum_{r \in \mathcal{S}^c} \|\bar{\boldsymbol{\xi}}^{\mathcal{R}} - \boldsymbol{\xi}^r\| \leq \sqrt{|\mathcal{S}^c|} \sqrt{\frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\bar{\boldsymbol{\xi}}^{\mathcal{R}} - \boldsymbol{\xi}^r\|^2}.$$

For interpretation, we denote $\nu^2 := \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\bar{\boldsymbol{\xi}}^{\mathcal{R}} - \boldsymbol{\xi}^r\|^2$, which can be interpreted as the variance in optimal dual variables of our problem. Then, the term on the right-hand side of the inequality above can be viewed as a bootstrap estimator of ν , which intuitively converges to ν as $\mathcal{S}^c \rightarrow \mathcal{R}$. To formalize this intuition, let us expand the squared norm term and apply the triangle inequality:

$$\begin{aligned} \left| \frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\bar{\boldsymbol{\xi}}^{\mathcal{R}} - \boldsymbol{\xi}^r\|^2 - \nu^2 \right| &= \left| \frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\boldsymbol{\xi}^r\|^2 - 2 \langle \bar{\boldsymbol{\xi}}^{\mathcal{R}}, \bar{\boldsymbol{\xi}}^{\mathcal{S}^c} \rangle - \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\boldsymbol{\xi}^r\|^2 + 2 \langle \bar{\boldsymbol{\xi}}^{\mathcal{R}}, \bar{\boldsymbol{\xi}}^{\mathcal{R}} \rangle \right| \\ &\leq \left| \frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\boldsymbol{\xi}^r\|^2 - \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\boldsymbol{\xi}^r\|^2 \right| + 2 \left| \langle \bar{\boldsymbol{\xi}}^{\mathcal{R}}, \bar{\boldsymbol{\xi}}^{\mathcal{R}} - \bar{\boldsymbol{\xi}}^{\mathcal{S}^c} \rangle \right| \\ &\leq \left| \frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\boldsymbol{\xi}^r\|^2 - \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\boldsymbol{\xi}^r\|^2 \right| + 2M \|\bar{\boldsymbol{\xi}}^{\mathcal{R}} - \bar{\boldsymbol{\xi}}^{\mathcal{S}^c}\| \end{aligned}$$

By Bardenet and Maillard (2015, corollary 2.5) again, there exists $M_2 > 0$ such that, with probability $1 - \delta$,

$$\left| \frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\boldsymbol{\xi}^r\|^2 - \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\boldsymbol{\xi}^r\|^2 \right| \leq M_2 \sqrt{\left(\frac{1}{|\mathcal{S}^c|} - \frac{1}{|\mathcal{R}|} \right) \log(1/\delta)},$$

and $\|\bar{\xi}^{\mathcal{R}} - \bar{\xi}^{\mathcal{S}^c}\|$ satisfies a similar inequality as (16). All together, with probability $1 - 2\delta$,

$$\left| \frac{1}{|\mathcal{S}^c|} \sum_{r \in \mathcal{S}^c} \|\bar{\xi}^{\mathcal{R}} - \xi^r\|^2 - \nu^2 \right| \leq M_2 \sqrt{\left(\frac{1}{|\mathcal{S}^c|} - \frac{1}{|\mathcal{R}|}\right) \log(1/\delta)} + M_1 \sqrt{d \left(\frac{1}{|\mathcal{S}^c|} - \frac{1}{|\mathcal{R}|}\right) \log(1/\delta)},$$

yielding

$$\sum_{r \in \mathcal{S}^c} \|\bar{\xi}^{\mathcal{R}} - \xi^r\| \leq \sqrt{|\mathcal{S}^c|} \nu + \sqrt{|\mathcal{S}^c|} M_3 \left(\frac{1}{|\mathcal{S}^c|} - \frac{1}{|\mathcal{R}|}\right)^{1/4} (\log(1/\delta))^{1/4}, \quad (17)$$

with $M_3 := M_2 + M_1 \sqrt{d}$.

Combining (16) and (17), we obtain that, with probability $1 - 3\delta$ over the sample \mathcal{S} ,

$$\sum_{r \in \mathcal{S}^c} \left| q(z, \bar{\xi}^{\mathcal{S}}, \mathbf{d}^r) - q(z, \xi^r, \mathbf{d}^r) \right| \leq L \sqrt{|\mathcal{S}^c|} \nu + E$$

where E is a bootstrap error term equal to

$$\begin{aligned} E &= L |\mathcal{S}^c| M_1 \sqrt{d \left(\frac{1}{|\mathcal{S}|} - \frac{1}{|\mathcal{R}|}\right) \log(1/\delta)} + \sqrt{|\mathcal{S}^c|} M_3 L \left(\frac{1}{|\mathcal{S}^c|} - \frac{1}{|\mathcal{R}|}\right)^{1/4} (\log(1/\delta))^{1/4} \\ &\leq \sqrt{|\mathcal{S}^c|} M_3 L \left[\sqrt{|\mathcal{R}|} \left(\frac{1}{|\mathcal{S}|} - \frac{1}{|\mathcal{R}|}\right)^{1/2} + \left(\frac{1}{|\mathcal{S}^c|} - \frac{1}{|\mathcal{R}|}\right)^{1/4} \right] \sqrt{\log(1/\delta)}, \end{aligned}$$

because $M_3 \geq M_1 \sqrt{d}$ and for δ such that $\log(1/\delta) > 1$.

To conclude the proof, let us observe that $M_3 = M_2 + M_1 \sqrt{d}$ and $d = 2|\mathcal{E}| + |\mathcal{N}| \times |\mathcal{K}|$. \square

Appendix B: Effect of Regularizer γ

The regularizing constant γ plays a crucial role in the performance of the algorithm. An appropriate value of γ is essential for achieving optimal convergence and solution quality. When set too high, the regularizing term has minimal impact on the objective function, making the problem more challenging to solve. Conversely, when set too low, the regularizing term dominates the objective function, resulting in easier but less accurate solutions.

To illustrate this, Figure B.1 presents the results of our experiments on **R** instances with 160 scenarios, where we vary the value of γ . It displays the average runtime and the objective value for each value of γ to provide insights on how the choice of γ affects the solution quality and computational performance of the algorithm. In the experimental results presented in Section 4, we selected an appropriate value for the regularizing constant γ in order to strike a balance between the two extremes of the spectrum, as depicted in Figure B.1.

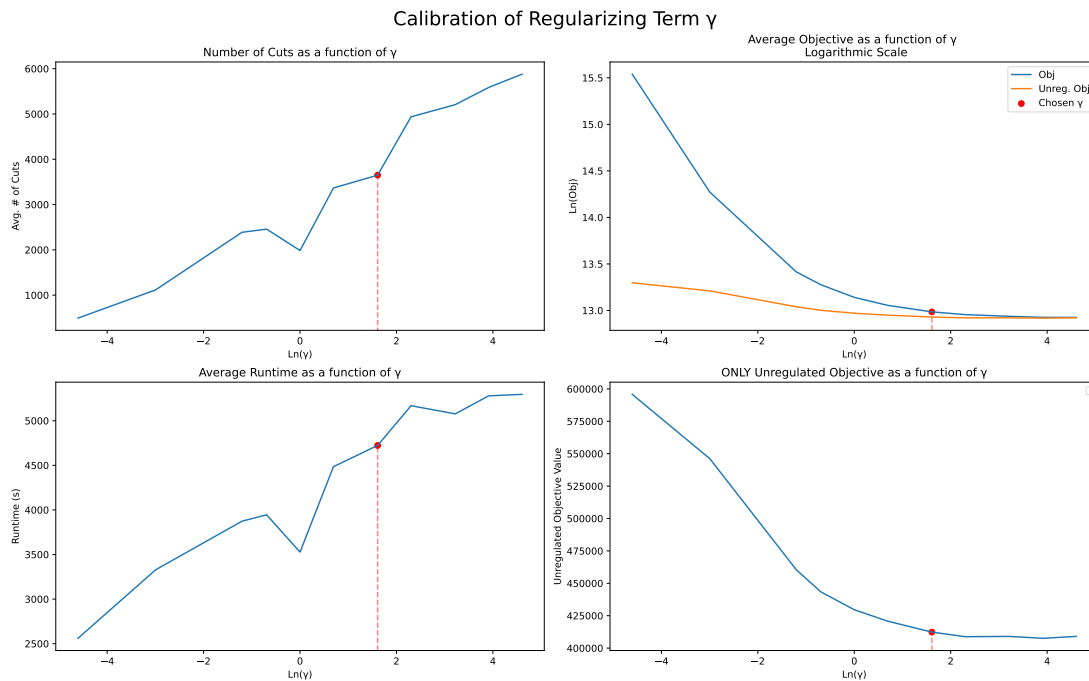


Figure B.1 Effect of Regularizer on the Algorithm's Performance.

Appendix C: Additional Numerical Results

In this section, we present additional numerical results that complement the results in Section 4.

C.1. Comparison of Different Stochastic Cutting-Plane Algorithms

In Section 4.2, we benchmark the performance of different variants of the stochastic cutting plane algorithm (namely the multi-, single-, and k -cut algorithms) with different warm-starting strategies at the root node. Recall that we terminate our algorithm after 3,600 seconds or as soon as it achieves an optimality gap of with confidence level $\alpha = 0.90$.

Accordingly, the average computational time reported in Table 2 are capped at 3,600 seconds whenever the algorithm does not converge within this time limit. To appreciate this censoring issue, Table C.1 presents the fraction of instances solved to ϵ -optimality for each combination of algorithm and warm-start strategy.

Finally, we explore the behavior of the stochastic k -cut algorithm when we vary the number of epigraph variables k . We vary the ratio $\frac{k}{|\mathcal{R}|}$ and report the average runtime and optimality gap in Tables C.2 and C.3 respectively. We observe a trade-off between computational time and optimality gap: A smaller number of clusters generally results in faster computations, while a larger value for k leads to a smaller optimality gap. To find a compromise between computation speed and optimality, we opt for $\frac{k}{|\mathcal{R}|} = 0.5$ in our experiments in Section 4.

Table C.1 Percentage (in %) of instances for which the algorithm converged within the time limit (3,600 seconds), for the multi-, single-, and k -cut stochastic cutting plane algorithms, with different warm-start strategies at the root node. Metrics are averaged across instances with the same number of nodes $|\mathcal{N}|$.

$ \mathcal{N} $	Multi-Cut			Single-Cut			k -Cut		
	None	Multi	Single	None	Multi	Single	None	Multi	Single
10	100.00	100.00	80.00	100.00	100.00	100.00	60.00	100.00	100.0
30	46.67	66.67	60.00	33.33	33.33	33.33	40.00	66.67	40.0
50	60.00	66.67	66.67	33.33	33.33	66.67	53.33	66.67	60.0
70	40.00	50.00	50.00	40.00	40.00	45.00	50.00	50.00	50.0
100	45.00	50.00	50.00	35.00	40.00	45.00	50.00	50.00	50.0
150	30.00	50.00	50.00	45.00	35.00	55.00	50.00	50.00	50.0
200	45.00	50.00	45.00	50.00	50.00	50.00	45.00	50.00	50.0

Table C.2 Runtime (in seconds) of the k -cut stochastic cutting plane algorithm, with different fraction of clusters per number of scenarios, $\frac{k}{|\mathcal{R}|}$. Metrics are averaged across instances with the same number of nodes $|\mathcal{N}|$.

$ \mathcal{N} $	k -cut with $\frac{k}{ \mathcal{R} }$ clusters											
	0.001	0.05	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
10	63.26	67.54	72.17	82.25	96.80	119.31	86.34	152.75	132.56	187.49	180.51	191.19
30	2868.71	2759.48	2800.97	2807.02	2747.69	2745.31	2820.14	3220.98	3104.16	3142.78	3135.66	3531.56
50	2545.08	2656.20	2693.95	2630.63	2691.93	2722.20	2783.76	3049.77	3131.54	3435.85	3420.01	3719.07
70	2885.32	2827.23	2710.32	2776.65	2950.30	3091.63	3035.84	3374.40	3572.29	3673.15	3802.27	3841.61
100	3298.48	3121.10	3041.21	3138.44	3270.07	3314.94	3534.80	3682.95	3701.85	3912.04	3995.43	3997.15
150	3693.06	3357.10	3282.66	3288.66	3465.72	3602.85	3687.58	3979.59	4096.43	4129.80	4192.95	4189.60
200	3783.36	3142.35	3240.84	3230.62	3475.07	3654.78	3784.40	3779.90	4009.51	4080.76	4086.29	4047.38

Table C.3 Gap (in %) of the k -cut stochastic cutting plane algorithm, with different fraction of clusters per number of scenarios, $\frac{k}{|\mathcal{R}|}$. Metrics are averaged across instances with the same number of nodes $|\mathcal{N}|$.

$ \mathcal{N} $	k -cut with $\frac{k}{ \mathcal{R} }$ clusters											
	0.001	0.05	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
10	0.88	0.64	0.33	0.25	0.83	1.03	0.66	0.34	0.38	0.56	0.25	0.45
30	8.03	7.34	7.93	7.91	7.85	6.95	6.77	6.66	6.13	5.78	5.79	5.49
50	5.20	5.75	4.80	4.42	4.16	3.74	3.70	3.24	2.95	3.79	2.66	2.58
70	5.45	5.58	4.79	5.43	4.31	3.87	3.56	3.98	2.97	2.76	2.59	3.56
100	6.33	5.73	5.82	5.15	5.52	3.87	3.08	2.52	2.70	2.36	2.17	4.17
150	5.99	8.46	7.12	4.67	5.45	5.13	3.54	3.95	3.97	5.02	2.08	2.96
200	8.70	3.35	6.18	4.05	2.82	3.62	2.72	5.32	3.16	6.75	5.48	6.25

C.2. Benchmarking Scalability on Synthetic Instances

To verify the correctness of our implementation, we use the smallest instances to verify that all methods terminate with the same optimal solution. To this end, Table C.4 reports the optimality gap (in %) and

final objective value for each algorithm, averaged over instances with the same number of nodes $|\mathcal{N}|$ and for which Gurobi converged to within 5% of optimality.

Table C.4 Optimality gap (in %) and final objective value for each algorithm, averaged over synthetic instances with the same number of nodes $|\mathcal{N}|$, where Gurobi converged to within 5% of optimality.

$ \mathcal{N} $	Gurobi with (1)		Deterministic		Stochastic	
	Gap	Objective	Gap	Objective	Gap	Objective
10	0.00	10,502.04	0.07	10,502.05	0.23	10,509.23
30	2.33	153,128.18	0.95	152,705.24	1.94	152,824.44
50	2.09	427,137.66	0.94	426,360.59	1.49	426,424.82
70	2.01	738,840.46	0.75	733,340.07	0.82	731,619.70

Figure C.1 illustrates the scalability of our single-cut method with respect to the number of scenarios and commodities and depicts the optimality gap achieved depending on the total number of nodes $|\mathcal{N}|$ (horizontal axis), the number of scenarios $|\mathcal{R}|$, and the number of commodities $|\mathcal{K}|$. We observe that the complexity of the problem (measured in terms of the final optimality gap) increases with all three problem dimensions, with the number of commodities and nodes having the most noticeable impact.

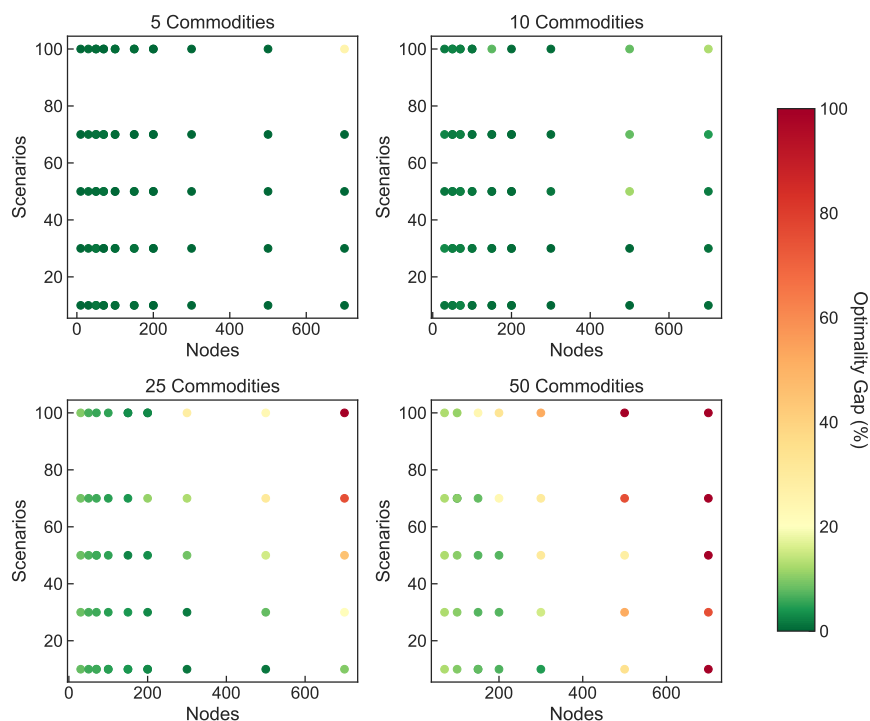


Figure C.1 Optimality gaps achieved by the single-cut stochastic cutting plane algorithm on all synthetic instances. For each combination of number of nodes $|\mathcal{N}|$, number of commodities $|\mathcal{K}|$, and number of scenarios $|\mathcal{R}|$, results are averaged across 3 random instances.

C.3. Benchmarking on the Instances from Crainic et al. (2000)

First, we provide further information on how we generate stochastic instances with a specified number of scenarios $|\mathcal{R}|$ from the \mathbf{R} instances of Crainic et al. (2000). As in Crainic et al. (2021b, 2016), Boland et al. (2016), we use the instances from classes 4 to 10 (9 instances per class), yielding 63 different instances (see Table 1 in Crainic et al. 2016, for a description of the instance classes). For each instance, we consider the same set of nodes, arcs, and commodities. We take the same arc capacities, edge construction costs, and flow transportation costs as those provided in the instances. For each commodity, the instance provides a vector of demand which we refer to as the nominal scenario and denote $\hat{\mathbf{b}}$. We generate scenarios by randomly inflating/deflating the nominal scenario, i.e., we sample $\mathbf{b} \sim \mathcal{U}(0.5, 1.2) \cdot \hat{\mathbf{b}}$. However, some of these demand scenarios might be infeasible because of arc capacity constraints. So, we only keep randomly generated scenarios that are feasible. We do so until we generate the desired number of scenarios, $|\mathcal{R}|$.