

Decomposition and Adaptive Sampling for Data-Driven Inverse Linear Optimization

Rishabh Gupta¹ and Qi Zhang ^{*1}

¹*Department of Chemical Engineering and Materials Science, University of Minnesota, Minneapolis, MN 55455, USA*

Abstract

This work addresses inverse linear optimization where the goal is to infer the unknown cost vector of a linear program. Specifically, we consider the data-driven setting in which the available data are noisy observations of optimal solutions that correspond to different instances of the linear program. We introduce a new formulation of the problem that, compared to other existing methods, allows the recovery of a less restrictive and generally more appropriate admissible set of cost estimates. It can be shown that this inverse optimization problem yields a finite number of solutions, and we develop an exact two-phase algorithm to determine all such solutions. Moreover, we propose an efficient decomposition algorithm to solve large instances of the problem. The algorithm extends naturally to an online learning environment where it can be used to provide quick updates of the cost estimate as new data becomes available over time. For the online setting, we further develop an effective adaptive sampling strategy that guides the selection of the next samples. The efficacy of the proposed methods is demonstrated in computational experiments involving two applications, customer preference learning and cost estimation for production planning. The results show significant reductions in computation and sampling efforts.

Keywords: inverse optimization, online learning, adaptive sampling

1 Introduction

Inverse optimization is an emerging new paradigm for uncovering hidden decision-making mechanisms from observed decision data. Following the principle of optimality (Schoemaker, 1991) that is commonly applied in various fields including economics, psychology, and evolutionary biology, the key idea in inverse optimization is to model a decision-making process as an optimization problem. Decisions can then be viewed as optimal or near-optimal solutions of an optimization model, and the inverse optimization problem (IOP) is to infer this model, if otherwise unknown, from observations. A major advantage of this approach is its ability to explicitly include constraints. This eases the incorporation of domain knowledge, which is often readily available in the form of constraints, significantly. As a result, compared to common black-box machine

*Corresponding author (qizh@umn.edu)

learning methods, inverse optimization offers the promise of models with enhanced prediction accuracy and interpretability.

The notion of inverse optimization was first introduced by [Burton & Toint \(1992\)](#) who considered the problem of determining travel costs on a network as perceived by the users given the routes they have taken. This has inspired research on several inverse network optimization problems ([Yang et al., 1997](#); [Zhang & Cai, 1998](#); [Zhang & Liu, 1999](#); [Liu & Zhang, 2006](#)). Since then, inverse optimization has found application in a myriad of fields, such as radiation therapy planning ([Chan et al., 2014](#); [Babier et al., 2018](#)), investment portfolio optimization ([Bertsimas et al., 2012](#)), electricity demand forecasting ([Saez-Gallego & Morales, 2018](#)), auction mechanism design ([Beil & Wein, 2003](#); [Birge et al., 2017](#)), biological systems ([Burgard & Maranas, 2003](#); [Arechavaleta et al., 2008](#); [Terekhov et al., 2010](#)), and optimal control ([Hempel et al., 2015](#); [Westermann et al., 2020](#)).

Early works in inverse optimization focus on determining an objective function that makes the observed decisions, given the constraints of the problem, exactly optimal. In their seminal paper, [Ahuja & Orlin \(2001\)](#) present a generalized solution method for inverse optimization with linear forward optimization problems (FOPs). Some of the later works extend the theory to consider conic ([Iyengar & Kang, 2005](#); [Zhang & Xu, 2010](#)), discrete ([Schaefer, 2009](#); [Wang, 2009](#); [Bulut & Ralphs, 2016](#)), and nonlinear ([Chow et al., 2014](#)) FOPs.

More recently, the research focus has shifted towards *data-driven* inverse optimization in which we observe an agent’s decisions in multiple instances, which can be viewed as instances of the same optimization problem that differ in their input parameter values ([Mohajerin Esfahani et al., 2018](#)). With data-driven inverse optimization, one has a much greater chance of learning an optimization model that has true predictive power with respect to future decisions in unseen instances. Here, the observations are generally considered to be noisy and the existing literature is limited to the case of convex FOPs. The main distinction among the various proposed formulations is in terms of the loss function employed to fit the data. Minimization of the slack required to make the noisy data satisfy an optimality condition is considered by [Keshavarz et al. \(2011\)](#), [Bertsimas et al. \(2015\)](#), and [Mohajerin Esfahani et al. \(2018\)](#). However, [Aswani et al. \(2018\)](#) show that this kind of loss function can lead to statistically inconsistent estimates and propose to minimize the sum of some norm of residuals with respect to the decision variables. In the data-driven context, most existing works assume that the observations are available in a single batch, while more recent contributions also address online learning environments in which the observations are made sequentially ([Bärmann et al., 2017](#); [Dong et al., 2018](#); [Shahmoradi & Lee, 2019](#)).

In this work, we consider data-driven inverse *linear* optimization with noisy observations in both batch and online learning settings. Here, the goal is to estimate the unknown cost vector of a linear program (LP). Inverse linear optimization constitutes an important class of IOPs as many decision-making problems can be formulated or approximated as LPs. Although inverse linear optimization falls into the broader category of inverse convex optimization for which established solution methods exist, these more general methods often yield overly restricted sets of admissible cost estimates when applied to inverse linear optimization with noisy data (as discussed in detail in Section 2). This limitation is to a great extent shared by tailored approaches specifically designed to solve the IOP in the linear case more efficiently ([Babier et al., 2018](#); [Chan et al., 2019](#)). We note that the majority of the inverse linear optimization literature does not consider the data-driven

case but focuses on the single-instance setting (with possibly multiple noisy observations).

Our proposed framework is designed to recover the complete set of admissible solutions for the inverse linear optimization problem, while incorporating the notion of a reference cost vector that represents the user’s prior belief, which facilitates the selection of an appropriate point estimate. Based on a polyhedral understanding of the problem, we propose a two-phase approach that separates the tasks of denoising the data and parameter estimation. To solve large instances of the IOP, we develop an exact decomposition algorithm, which processes the data sequentially and can hence also serve as an efficient update method in online inverse optimization. For the online setting, we further develop an adaptive sampling strategy that guides the selection of the next samples in an effort to reduce the amount of required data. While adaptive sampling is quite a mainstream idea in machine learning (Domingo et al., 2002; Chang et al., 2005; Cozad et al., 2014), to the best of our knowledge, it has not yet been considered in inverse optimization. We believe that the development of such a framework can go a long way in increasing the acceptance of inverse optimization as an alternative to black-box modeling methods, especially in situations where data acquisition is expensive or time-intensive.

The main contributions of this work are as follows:

1. We introduce a new general formulation of the data-driven inverse linear optimization problem that considers multiple noisy observations collected for multiple experiments, which are problem instances with different input parameter values. We highlight several geometrical properties of the problem, and show that by assuming that optimal solutions lie at the vertices of the feasible region, we can recover the complete set of admissible cost estimates.
2. We show that the proposed IOP formulation yields a finite number of solutions. We introduce a two-phase algorithm that can recover all such solutions. Furthermore, we show that under a very mild condition, the IOP is guaranteed to have a unique solution.
3. We develop an efficient sequential decomposition algorithm to solve large instances of the IOP. The algorithm directly extends itself to online inverse optimization where it can be used to provide quick updates of the cost estimate as new data becomes available.
4. We propose an effective adaptive sampling strategy that guides the choice of the experiments in online inverse optimization. The adaptive sampling problem is formulated as a mixed-integer nonlinear program (MINLP) for which we provide an efficient heuristic solution algorithm.
5. We demonstrate the effectiveness of the proposed framework through a comprehensive set of computational experiments, addressing the problems of customer preference learning and cost estimation for production planning. The results indicate that generally, reasonable prediction accuracies can be achieved with relatively small numbers of experiments. Also, one can observe significant reductions in solution time and required number of samples due to the proposed decomposition and adaptive sampling methods, respectively.

The remainder of this paper is organized as follows. In Section 2, we present a formal description of the inverse linear optimization problem. In Section 3, we propose a new formulation

that utilizes a reference cost vector to find reasonable cost estimates, discuss its properties, and develop a two-phase solution algorithm. Section 4 introduces an exact decomposition algorithm that allows the efficient solution of large instances of the IOP and naturally extends to online inverse optimization. Our proposed adaptive sampling framework is detailed in Section 5. In Section 6, results from the computational studies are presented. Finally, we conclude in Section 7.

2 Background of Inverse Linear Optimization

Consider a decision-making problem that can be represented as an LP of the following form:

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && c^\top x \\ & \text{subject to} && Ax \leq b, \end{aligned} \tag{FOP}$$

which we call the forward optimization problem (FOP). The cost vector $c \in \mathbb{R}^n$ is unknown; however, experiments perturbing values in $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$, subject to restrictions $(A, b) \in \Pi$, can be designed to help improve our estimate of c . The results of these experiments, which are assumed to be optimal solutions to (FOP), are observed with some random noise. For a specific experiment, multiple samples can be collected such that one can recover the true optimal solution with some confidence, i.e. obtain a denoised estimate. We assume that for any $(A, b) \in \Pi$, the polyhedron represented by $Ax \leq b$ is compact and nonempty. This is a mild assumption as in essentially all real-world problems, the decision variables are bounded.

Given observations, the inverse optimization problem (IOP) is to obtain an estimate of c , \hat{c} , such that the difference between the observations and the solutions obtained from solving (FOP) with \hat{c} as the cost vector is minimized. The IOP is commonly formulated as follows:

$$\begin{aligned} & \underset{\hat{c} \in \mathbb{R}^n, \hat{x}}{\text{minimize}} && \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_{ij}\| \\ & \text{subject to} && \hat{x}_{ij} \in \arg \min_{\tilde{x} \in \mathbb{R}^n} \left\{ \hat{c}^\top \tilde{x} : A_i \tilde{x} \leq b_i \right\} \quad \forall i \in \mathcal{I}, j \in \mathcal{J}_i, \end{aligned} \tag{1}$$

where \mathcal{I} is the set of experiments, where each experiment i is associated with inputs (A_i, b_i) , \mathcal{J}_i denotes the set of noisy observations for experiment i , and x_{ij} is the observed output for $j \in \mathcal{J}_i$. The objective is to choose \hat{c} and \hat{x} such that the loss function, which is defined as the sum of some norm of the residuals, is minimized. Formulation (1) generalizes existing variants of the IOP from the literature. Some consider a setting in which A and b cannot be changed, which leads to the case of $|\mathcal{I}| = 1$ (Chan & Lee, 2018; Chan et al., 2019). Others consider random sampling of A and b without assigning the samples to distinct sets corresponding to specific inputs (Aswani et al., 2018); in this case, we have $|\mathcal{I}| = N$ with N being the total number of samples, and $|\mathcal{J}_i| = 1$ for all $i \in \mathcal{I}$. In fact, splitting the set of samples into input-specific subsets does not affect a formulation like (1); however, it will be an essential feature of our proposed alternative approach (more in Section 3).

Problem (1) is a bilevel optimization problem and is typically solved by replacing its lower-level problems with their optimality conditions. While most existing works make use of strong duality (Aswani et al., 2018; Chan et al., 2019; Shahmoradi & Lee, 2019), some have also applied

reformulations based on the Karush-Kuhn-Tucker (KKT) conditions (Keshavarz et al., 2011; Saez-Gallego et al., 2016). In case of LPs, both optimality conditions are equivalent. In this work, we use a KKT-based approach as the duality-based formulation is nonlinear and nonconvex in the constraints due to the presence of bilinear terms, whereas the constraints of the KKT-based formulation can be linearized by introducing binary variables, which is advantageous from a computational standpoint. Thus, we arrive at the following mixed-integer reformulation of (1):

$$\underset{\hat{c} \in \mathbb{R}^n, \hat{x}, s, \lambda, z}{\text{minimize}} \quad \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_{ij}\| \quad (2a)$$

$$\text{subject to} \quad \hat{c} + A_i^\top \lambda_{ij} = 0 \quad \forall i \in \mathcal{I}, j \in \mathcal{J}_i \quad (2b)$$

$$A_i \hat{x}_{ij} + s_{ij} = b_i \quad \forall i \in \mathcal{I}, j \in \mathcal{J}_i \quad (2c)$$

$$\lambda_{ij} \leq M z_{ij} \quad \forall i \in \mathcal{I}, j \in \mathcal{J}_i \quad (2d)$$

$$s_{ij} \leq M(e - z_{ij}) \quad \forall i \in \mathcal{I}, j \in \mathcal{J}_i \quad (2e)$$

$$\hat{x}_{ij} \in \mathbb{R}^n, s_{ij} \in \mathbb{R}_+^m, \lambda_{ij} \in \mathbb{R}_+^m, z_{ij} \in \{0, 1\}^m \quad \forall i \in \mathcal{I}, j \in \mathcal{J}_i, \quad (2f)$$

where M is a sufficiently large parameter and e denotes the all-ones vector. Constraints (2b), (2c), and (2d)–(2e) correspond to the stationarity, primal feasibility, and complementary slackness conditions, respectively. The following theorem characterizes the solution set of (2).

Theorem 1. *For a given feasible \hat{x} , the set of feasible \hat{c} in problem (2) is a polyhedral cone.*

Proof. Consider (2) for a specific experiment i and sample $j \in \mathcal{J}_i$ and a corresponding feasible \hat{x}_{ij} . Let \mathcal{K}_i be the set of constraints for experiment i , i.e. resulting from A_i and b_i . From (2f), we have that $\lambda_{ij} \geq 0$; hence, if $\lambda_{ijk} = 0$ for all $k \in \mathcal{K}_i$, then (2b) imply that $\hat{c} = 0$. Otherwise, if $\exists k \in \mathcal{K}_i$ such that $\lambda_{ijk} > 0$, then from constraints (2d)–(2e), we have that $s_{ijk} = 0$, and (2c) imply that \hat{x}_{ij} is such that $a_{ik}^\top \hat{x}_{ij} = b_{ik}$, where $a_{ik} \in \mathbb{R}^n$ defines the k th row of A_i . Hence, from (2b), we have that $\hat{c} \in \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_{ij})})$, where $\mathcal{T}(\hat{x}_{ij})$ denotes the set of constraints active at \hat{x}_{ij} . Since $\mathcal{T}(\hat{x}_{ij})$ is a finite set, the feasible region for \hat{c} associated with experiment i and sample j is a polyhedral cone. As this statement holds for every $i \in \mathcal{I}$ and $j \in \mathcal{J}_i$, we have

$$\hat{c} \in \bigcap_{i \in \mathcal{I}, j \in \mathcal{J}_i} \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_{ij})}),$$

which is the intersection of a finite number of polyhedral cones, hence also a polyhedral cone. \square

Admissible Set As a consequence of Theorem 1, there is no unique solution to problem (1), or equivalently (2), since for any optimal \hat{c} , $\alpha \hat{c}$ with α being any nonnegative scalar different than 1 yields another optimal solution. Instead, by means of Theorem 1, we can determine the full set of “inverse-optimal” \hat{c} , which we refer to as the *admissible set*. Also, to avoid the trivial solution $\hat{c} = 0$, it is common to use a slight variant of (1) that restricts the length of \hat{c} by adding a norm constraint, e.g. $\|\hat{c}\|_p = 1$ (Mohajerin Esfahani et al., 2018; Chan et al., 2019; Shahmoradi & Lee, 2019). However, the use of a norm constraint introduces additional nonconvexity into the problem formulation. While the choice of p -norm has been arbitrary, under some special conditions on c , the 1-norm and ∞ -norm have been shown to lead to tractable formulations (Chan et al., 2019).

Noisy Observations As (FOP) is an LP, with a nonzero c , any optimal solution will lie on the boundary of the polyhedral feasible region. Problem (1) with a p -norm constraint on c can be interpreted as the projection of noisy observations onto one of the polyhedron’s facets such that the total projection distance is minimized (Chan et al., 2019). While this approach provides good solutions when the FOP is strongly convex, it often leads to a severely restricted admissible set when the FOP is an LP. As illustrated in Figure 1, even if the true solution lies at a vertex, noise in the data can cause them to get projected onto one of the facets, making a vector orthogonal to that facet the only feasible \hat{c} . As highlighted by Shahmoradi & Lee (2019), this formulation also leads to unstable predictions in the presence of outliers in the data.

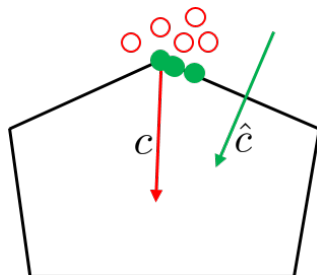


Figure 1: The polytope represents the feasible region of the FOP defined by $A_i x \leq b_i$, noisy data samples x_{ij} are depicted by red hollow circles, and their denoised estimates \hat{x}_{ij} are shown as green filled circles. The true cost vector is c , whereas \hat{c} is the solution to (1) with an additional p -norm constraint on \hat{c} .

Reference Cost It is important to note that the problem of estimating c given noisy observations consists of two tasks that have to be performed simultaneously: (i) denoising the data to obtain estimates of the real optimal solutions, and (ii) using these solutions to estimate c . Formulations of the form (1) and existing variants thereof mainly address the first aspect, but have deficiencies, as elucidated above, when it comes to finding a good point estimate of c , especially when the admissible set is large or the number of observations is small. Traditionally, when sampling is deterministic, inverse optimization has been facilitated by using a reference cost vector \bar{c} and employing an objective function of the form $\|\bar{c} - \hat{c}\|_2^2$ (Ahuja & Orlin, 2001; Heuberger, 2004). Such a reference cost represents a prior belief that is available in most practical applications, e.g. obtained through first-principles modeling or expert knowledge. A cost estimate \hat{c} that is close to a reference \bar{c} is often desired. It seems that the notion of such a reference, which can aid the process of recovering the real cost vector, has largely been ignored in the noisy case.

3 A Two-Phase Approach to Inverse Optimization

In this section, we propose a general inverse linear optimization model that finds a nontrivial estimate of the cost vector that, among all the ones that minimize the loss function, most closely resembles a reference cost vector. We discuss the main properties of this IOP and develop an exact two-phase solution algorithm.

3.1 Problem Formulation and Properties

We first consider the following problem:

$$\begin{aligned} & \underset{\hat{c}}{\text{minimize}} && \|\bar{c} - \hat{c}\|_2^2 \\ & \text{subject to} && \hat{c} \in \mathcal{C}, \end{aligned} \quad (3)$$

where \bar{c} is the reference cost vector, $\mathcal{C} = \mathcal{C}' \cup \{0\}$, and \mathcal{C}' is defined as follows:

$$\mathcal{C}' := \left\{ \hat{c} : (\hat{c}, \hat{x}) \in \underset{\hat{c} \in \mathbb{R}^n \setminus \{0\}, \hat{x}}{\arg \min} \left\{ \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_{ij}\| : \hat{x}_{ij} \in \underset{\tilde{x} \in \mathbb{R}^n}{\arg \min} \{ \hat{c}^\top \tilde{x} : A_i \tilde{x} \leq b_i \} \forall i \in \mathcal{I}, j \in \mathcal{J}_i \right\} \right\}. \quad (4)$$

We assume that for all $i \in \mathcal{I}$, (A_i, b_i) is chosen from some set Π such that the polyhedron represented by $A_i x \leq b_i$ is compact and nonempty.

The set \mathcal{C}' consists of all \hat{c} that are optimal in (1) except for the trivial solution $\hat{c} = 0$. However, instead of \mathcal{C}' , problem (3) considers \mathcal{C} as the set of admissible cost estimates, which does include the trivial solution. The rationale behind this setup is the following: When determining \mathcal{C}' , we do not want to consider $\hat{c} = 0$, which allows a minimum loss but does not provide any useful information. However, admitting the trivial solution in (3) can be helpful since the unlikely case in which $\hat{c} = 0$ is the optimal solution to (3) would immediately indicate that \bar{c} is a very bad reference. In addition, as discussed later in this section, admitting \mathcal{C} in (3) results in useful theoretical properties.

As mentioned in Section 2, solving problem (3) provides good results in the deterministic case; however, it fails when noisy data are considered since projection of data onto a facet causes \mathcal{C} to be a single ray. To overcome this issue, we propose to consider, instead of \mathcal{C}' , the following slightly modified set:

$$\hat{\mathcal{C}}' := \left\{ \hat{c} : (\hat{c}, \hat{x}) \in \underset{\hat{c} \in \mathbb{R}^n \setminus \{0\}, \hat{x}}{\arg \min} \left\{ \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_{ij}\| : \hat{x}_{ij} \in \underset{\tilde{x} \in \mathbb{R}^n}{\arg \min} \{ \hat{c}^\top \tilde{x} : A_i \tilde{x} \leq b_i \} \cap \mathcal{V}_i \forall i \in \mathcal{I} \right\} \right\}, \quad (5)$$

where \mathcal{V}_i denotes the set of extreme points of $\{\tilde{x} : A_i \tilde{x} \leq b_i\}$. This leads to the following IOP:

$$\begin{aligned} & \underset{\hat{c}}{\text{minimize}} && \|\bar{c} - \hat{c}\|_2^2 \\ & \text{subject to} && \hat{c} \in \hat{\mathcal{C}}, \end{aligned} \quad (\text{IOP})$$

where $\hat{\mathcal{C}} = \hat{\mathcal{C}}' \cup \{0\}$. The set $\hat{\mathcal{C}}'$ considers the projection of the data for each experiment i to one of the vertices of $\{\tilde{x} : A_i \tilde{x} \leq b_i\}$ such that the total projection distance is minimized. We consider two different cases depending on the nature of (FOP) to argue why this approach leads to a more appropriate admissible set. For the ease of exposition, our discussion is restricted to the case of a single experiment, i.e. $|\mathcal{I}| = 1$, with multiple samples.

1. If (FOP) has a unique optimal solution, the noisy samples are likely to be located close to the vertex at which the optimal solution lies (see Figure 2a). By solving the loss minimization problem in (5), we obtain a vertex \hat{x}_i , and the resulting $\hat{\mathcal{C}}$ is an exhaustive set of cost vectors that can make \hat{x}_i optimal for (FOP). Recall from Theorem 1 that any $\hat{c} \in \hat{\mathcal{C}}$ can be expressed as a conic combination of the vectors orthogonal to the facets associated with the constraints

active at \hat{x}_i . Thus, if \hat{x}_i is the “right” vertex, i.e. the true optimal solution, the true cost vector c will be one of the vectors in $\hat{\mathcal{C}}$. This is in contrast to \mathcal{C} , which in this case would be a single ray that does not contain c .

2. If (FOP) has multiple optimal solutions, the noisy samples are likely to be located close to the facet that represents the set of optimal solutions (see Figure 2b). If solving the loss minimization problem in (5) results in a vertex \hat{x}_i that is an optimal solution to (FOP), c will be one of the extreme rays of $\hat{\mathcal{C}}$. For this case, in most practical instances, we expect the data to show a preference toward a particular vertex of the facet, i.e. there is a unique optimal \hat{x}_i . Even if multiple \hat{x}_i achieve the same loss, we still have $c \in \hat{\mathcal{C}}$ as long as at least one of the \hat{x}_i is in fact an optimal solution to (FOP). In this case, one may argue that $\hat{\mathcal{C}}$ is overly large compared to \mathcal{C} ; however, if $\bar{c} = c$, the true cost vector will be recovered when solving (IOP).

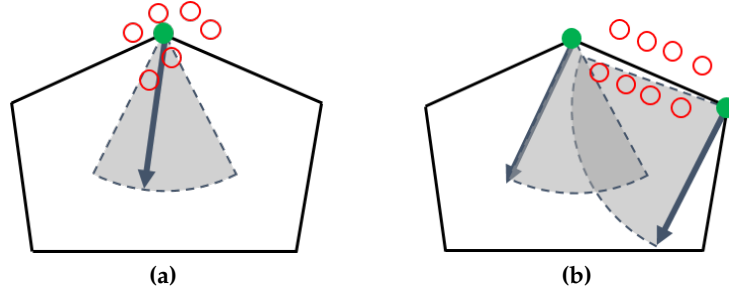


Figure 2: Different cases when projection of noisy data is restricted to one of the vertices of the polyhedral feasible region. Arrows show the true cost vector c , noisy data samples x_{ij} are depicted by red hollow circles, and their denoised estimates obtained from solving the loss minimization problem in (5) are shown as green filled circles. The gray shaded regions indicate the set $\hat{\mathcal{C}}$.

Lemma 1. *The set $\hat{\mathcal{C}}$ is nonempty.*

Proof. Since $\{\tilde{x} : A_i \tilde{x} \leq b_i\}$ is assumed to be compact and nonempty, $\arg \min_{\tilde{x} \in \mathbb{R}^n} \{\hat{c}^\top \tilde{x} : A_i \tilde{x} \leq b_i\}$ is nonempty, and for every $\hat{c} \in \mathbb{R}^n$, there is at least one vertex of $\{\tilde{x} : A_i \tilde{x} \leq b_i\}$ that minimizes $\hat{c}^\top \tilde{x}$. This implies that $\arg \min_{\tilde{x} \in \mathbb{R}^n} \{\hat{c}^\top \tilde{x} : A_i \tilde{x} \leq b_i\} \cap \mathcal{V}_i$ is nonempty and finite, due to the finiteness of \mathcal{V}_i , for all $i \in \mathcal{I}$ and $\hat{c} \in \mathbb{R}^n$. As a result, there is at least one (\hat{c}, \hat{x}) that minimizes the loss function in the definition of the outer arg min set in (5); hence, $\hat{\mathcal{C}}$ is nonempty. \square

Theorem 2. *The set of optimal solutions to (IOP) is finite.*

Proof. We start by characterising the set $\hat{\mathcal{C}}$. Consider the minimization problem that describes $\hat{\mathcal{C}}$:

$$\begin{aligned} & \underset{\hat{c} \in \mathbb{R}^n \setminus \{0\}, \hat{x}}{\text{minimize}} && \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_i\| \\ & \text{subject to} && \hat{x}_i \in \arg \min_{\tilde{x} \in \mathbb{R}^n} \left\{ \hat{c}^\top \tilde{x} : A_i \tilde{x} \leq b_i \right\} \cap \mathcal{V}_i \quad \forall i \in \mathcal{I}. \end{aligned} \tag{6}$$

As stated in Lemma 1, the solution set of problem (6) is nonempty. Furthermore, as \mathcal{V}_i is finite and $\hat{x}_i \in \mathcal{V}_i$, the set of feasible \hat{x}_i is finite. Since we also have a finite number of experiments, the set

of feasible \hat{x} to (6) is finite. Hence, the set of optimal \hat{x} , which we denote by $\hat{\mathcal{X}}^*$, is finite. From Theorem 1, the set of feasible \hat{c} for any $\hat{x} \in \hat{\mathcal{X}}^*$ is $\bigcap_{i \in \mathcal{I}} \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i)}) \setminus \{0\}$, which is also the set of optimal \hat{c} for that given optimal \hat{x} as \hat{c} does not appear in the objective function of (6). It follows that, considering all $\hat{x} \in \hat{\mathcal{X}}^*$, the set $\hat{\mathcal{C}}$ can be expressed as follows:

$$\hat{\mathcal{C}} = \bigcup_{\hat{x} \in \hat{\mathcal{X}}^*} \left\{ \bigcap_{i \in \mathcal{I}} \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i)}) \right\} \setminus \{0\}. \quad (7)$$

Problem (IOP) can then be solved by solving the following problem for every $\hat{x} \in \hat{\mathcal{X}}^*$:

$$\begin{aligned} & \underset{\hat{c}}{\text{minimize}} && \|\bar{c} - \hat{c}\|_2^2 \\ & \text{subject to} && \hat{c} \in \bigcap_{i \in \mathcal{I}} \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i)}), \end{aligned} \quad (8)$$

which is the minimization of a strictly convex function over a convex feasible region and hence has a unique optimal solution, which we denote by $\tilde{c}(\hat{x})$. The set $\tilde{\mathcal{C}} := \{\tilde{c}(\hat{x})\}_{\hat{x} \in \hat{\mathcal{X}}^*}$ is finite. Hence, the set of optimal solutions to (IOP), which can be expressed as

$$\tilde{\mathcal{C}}^* := \left\{ \hat{c} : \hat{c} \in \tilde{\mathcal{C}}, \|\bar{c} - \hat{c}\|_2^2 = \min_{\hat{x} \in \hat{\mathcal{X}}^*} \|\bar{c} - \tilde{c}(\hat{x})\|_2^2 \right\}, \quad (9)$$

is also finite. □

Condition 1. There is a unique optimal \hat{x} to problem (6).

At first glance, Condition 1 seems to be very restrictive. But in fact, it holds in almost all practical instances. Consider an instance in which at any optimal solution to (6), there is a unique optimal \hat{x}_i for each experiment i except for one experiment p . Let \hat{x}_i^* be the unique optimal \hat{x}_i for all $i \in \mathcal{I} \setminus \{p\}$, and let $\hat{\mathcal{X}}_p^*$ denote the set of multiple optimal \hat{x}_p . Then, the following two conditions have to hold: (i) The resulting loss associated with experiment p , i.e. $\sum_{j \in \mathcal{J}_p} \|x_{pj} - \hat{x}_p\|_2^2$, is the same for all $\hat{x}_p \in \hat{\mathcal{X}}_p^*$. (ii) For every $\hat{x}_p \in \hat{\mathcal{X}}_p^*$, there exists a \hat{c} that renders all \hat{x}_i^* , $i \in \mathcal{I} \setminus \{p\}$, and \hat{x}_p optimal for the FOP, which is equivalent to the following condition:

$$\text{cone}(\{-a_{pt}\}_{t \in \mathcal{T}(\hat{x}_p)}) \cap \left(\bigcap_{i \in \mathcal{I} \setminus \{p\}} \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i^*)}) \right) \neq \{0\} \quad \forall \hat{x}_p \in \hat{\mathcal{X}}_p^*. \quad (10)$$

While already the first condition is very unlikely to hold in practice, the second also becomes more improbable for $|\hat{\mathcal{X}}_p^*| > 1$ as the number of experiments increases. Hence, we conclude that the case of $|\hat{\mathcal{X}}^*| > 1$ is highly unlikely, which is consistent with our observation in our computational experiments. In theory, however, especially if $|\mathcal{I}|$ is small, there is the possibility that Condition 1 does not hold. In the next subsection, we present a two-phase algorithm that is guaranteed to find the complete set of optimal solutions to (IOP) even if Condition 1 does not hold.

Corollary 1. If Condition 1 holds, (IOP) has a unique solution.

Proof. Condition 1 implies that $|\hat{\mathcal{X}}^*| = 1$. Corollary 1 then directly follows from Theorem 2. □

Thus, in most cases, solving (IOP) will result in a unique \hat{c} that is closest to the prior belief \bar{c} . This is in contrast to other approaches in the literature where a variant of (6) is solved that results in a set of solutions from which a \hat{c} is randomly selected. Note that while all $\hat{c} \in \hat{\mathcal{C}}$ achieve the same prediction accuracy on the training set, they may not show the same performance on unseen data. Our approach of incorporating a prior belief \bar{c} resolves this ambiguity and determines whether \bar{c} is in the admissible set or, if not, how much it is outside the admissible set. This is a generally desirable feature in practice.

3.2 Two-Phase Algorithm

We start by presenting tractable reformulations of (6) and (8), and then show how they can be combined to obtain the set of optimal solutions to (IOP).

Applying a KKT-based approach, we obtain the following reformulation of (6):

$$\begin{aligned} & \underset{\hat{c}, \hat{x}, s, \lambda, z, w, \hat{c}^+, \hat{c}^-}{\text{minimize}} && \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_i\| && \text{(P1a)} \end{aligned}$$

$$\text{subject to} \quad \hat{c} + A_i^\top \lambda_i = 0 \quad \forall i \in \mathcal{I} \quad \text{(P1b)}$$

$$A_i \hat{x}_i + s_i = b_i \quad \forall i \in \mathcal{I} \quad \text{(P1c)}$$

$$\lambda_i \leq M z_i \quad \forall i \in \mathcal{I} \quad \text{(P1d)}$$

$$s_i \leq M(e - z_i) \quad \forall i \in \mathcal{I} \quad \text{(P1e)}$$

$$e^\top z_i \geq n \quad \forall i \in \mathcal{I} \quad \text{(P1f)}$$

$$\hat{c} = \hat{c}^+ - \hat{c}^- \quad \text{(P1g)}$$

$$\hat{c}^+ \leq w \quad \text{(P1h)}$$

$$\hat{c}^- \leq e - w \quad \text{(P1i)}$$

$$e^\top (\hat{c}^+ + \hat{c}^-) = 1 \quad \text{(P1j)}$$

$$\hat{x}_i \in \mathbb{R}^n, s_i \in \mathbb{R}_+^m, \lambda_i \in \mathbb{R}_+^m, z_i \in \{0, 1\}^m \quad \forall i \in \mathcal{I} \quad \text{(P1k)}$$

$$\hat{c} \in \mathbb{R}^n, \hat{c}^+ \in \mathbb{R}_+^n, \hat{c}^- \in \mathbb{R}_+^n, w \in \{0, 1\}^n, \quad \text{(P1l)}$$

where M is a sufficiently large parameter. Constraints (P1b)–(P1e) correspond to the KKT optimality conditions of the lower-level problems in (6), (P1f) ensures that \hat{x}_i is a vertex of the polyhedron $\{\tilde{x} : A_i \tilde{x} \leq b_i\}$, and (P1g)–(P1j) represent a linearization of the condition $\|\hat{c}\|_1 = 1$, which excludes $\hat{c} = 0$ from the set of feasible solutions. Note that (P1) is not an exact reformulation of (6) since the constraint $\|\hat{c}\|_1 = 1$ cuts off more than just the single point $\hat{c} = 0$. However, the following property suggests that this restriction does not affect the solution set of (IOP).

Lemma 2. *Problems (6) and (P1) have the same set of feasible \hat{x} .*

Proof. By construction, for a given feasible \hat{c} , (6) and (P1) exhibit the same set of feasible \hat{x} . Observe that any \hat{c} for which $\|\hat{c}\|_1 \neq 1$ is infeasible in (P1). Let \tilde{c} denote a \hat{c} that is feasible in (6) but not in (P1), i.e. $\|\tilde{c}\|_1 \neq 1$. However, there exists a positive α such that $\hat{c} = \alpha \tilde{c}$ is feasible in (P1), and it exhibits the same set of feasible \hat{x} as \tilde{c} does in (6) since uniform positive scaling of a cost vector does not change the set of optimal solutions of an LP. Hence, problems (6) and (P1) have the same set of feasible \hat{x} . \square

An exact reformulation of (8) directly follows from the definition of polyhedral cones:

$$\begin{aligned}
& \underset{\hat{c}, \gamma}{\text{minimize}} && \|\bar{c} - \hat{c}\|_2^2 \\
& \text{subject to} && \hat{c} = - \sum_{t \in \mathcal{T}(\hat{x}_i)} \gamma_{it} a_{it}^\top \quad \forall i \in \mathcal{I} \\
& && \gamma_{it} \geq 0 \quad \forall i \in \mathcal{I}, t \in \mathcal{T}(\hat{x}_i),
\end{aligned} \tag{P2}$$

where the constraints require that \hat{c} can be expressed as a conic combination of $-a_{it}, t \in \mathcal{T}(\hat{x}_i)$, for every $i \in \mathcal{I}$.

Proposition 1. *If Condition 1 holds, the optimal solution to (IOP) can be obtained by solving (P1), which provides the unique optimal \hat{x}^* , and subsequently solving (P2) with $\hat{x} = \hat{x}^*$.*

Proof. Proposition 1 directly follows from Corollary 1 and Lemma 2. \square

In the general case where Condition 1 does not necessarily hold (or where we do not know in advance that it holds), the complete set $\hat{\mathcal{X}}^*$ can be recovered by re-solving (P1) with added integer cuts of the following form until the objective function value changes:

$$\sum_{i \in \mathcal{I}} \sum_{t \in \mathcal{T}(\hat{x}_i)} z_{it} \leq n|\mathcal{I}| - 1. \tag{11}$$

Since the set of active constraints at a particular \hat{x}_i is given by the values of the binary variables z_i , we use (11) to impose the condition that for a different optimal solution to exist, at least one of the \hat{x}_i has to induce a different set of active constraints. A two-phase algorithm incorporating integer cuts to obtain the complete set of optimal solutions to (IOP) as described in the proof of Theorem 2 is shown in Algorithm 1.

Algorithm 1 Two-phase algorithm for solving (IOP).

PHASE 1:

- 1: solve (P1), obtain $\hat{x}^*, r^* \leftarrow \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_i^*\|$
- 2: initialize: $k \leftarrow 1, r^1 \leftarrow r^*, \hat{x}^1 \leftarrow \hat{x}^*$
- 3: **while** $r^k = r^*$ **do**
- 4: add integer cut (11) for \hat{x}^k to (P1)
- 5: solve (P1), obtain $\hat{x}^*, \hat{x}^{k+1} \leftarrow \hat{x}^*, r^{k+1} \leftarrow \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_i^{k+1}\|$
- 6: $k \leftarrow k + 1$
- 7: **end while**

PHASE 2:

- 8: $K \leftarrow k - 1$
 - 9: **for all** $k = 1, \dots, K$ **do**
 - 10: solve (P2) with $\hat{x} = \hat{x}^k$, obtain $\hat{c}^*, \hat{c}^k \leftarrow \hat{c}^*, v^k \leftarrow \|\bar{c} - \hat{c}^k\|_2^2$
 - 11: **end for**
 - 12: **return** all \hat{c}^k for which $v^k = \min_{k'=1, \dots, K} v^{k'}$
-

Remark 1. We compare our approach with that of [Shahmoradi & Lee \(2019\)](#) where data is denoised by choosing the projection that maximizes the number of active constraints. In their proposed method, outliers are removed from the data by considering samples within a specified

quantile statistic (θ) of the optimality error. Their method expects some domain knowledge to estimate the expected level of noise and assign an appropriate value of θ and total optimality error (E). As demonstrated in their work, if the hyperparameters θ and E are not chosen properly, the solution set may exclude the true cost vector c . In contrast, our approach, by assuming that an optimal solution to (FOP) is one of the vertices of the feasible region, eliminates the need of any hyperparameters and considers the most conservative (largest) \widehat{C} to choose a \hat{c} from. It relies on a reference cost vector \bar{c} , which reflects the user's prior belief, to act as an anchor to effectively navigate the larger search space and recover a good cost estimate.

Remark 2. Problems (P1) and (P2) do not include any additional constraints on the values of c (and therefore \hat{c}) apart from it being a nonzero vector in (P1). In practice, it is likely that we have additional information about the true cost vector. For example, if the cost coefficients to be estimated represent costs of items, one can safely assume them to be strictly positive. Such information can be directly incorporated into (P1) and (P2) in the form of additional constraints, which can further restrict the admissible set and help obtain a reasonable cost estimate. This can be especially useful in situations where the observed data has a high variance or the vertices of the polyhedron are very close to each other. We illustrate this point through an example (see Figure 3): Let the (FOP) be $\max_{x_1, x_2} \{2x_1 + x_2 : 14x_1 + 10x_2 \leq 17, 0 \leq x_1, x_2 \leq 1\}$. Let $|\mathcal{I}| = 1$ and $|\mathcal{J}| = 4$ with $x_1 = [(1.05, 0.26), (0.9, 0.15), (1.2, 0.15), (1, 0.03)]$. Clearly, the true solution is $(1, 0.3)$. However, (P1) yields the denoised estimate $\hat{x}_1 = (1, 0)$, whereas, if we consider additional constraints stating $\hat{c} < 0$ (note that the FOP here is a maximization problem), the estimate becomes $\hat{x}_1 = (1, 0.3)$.

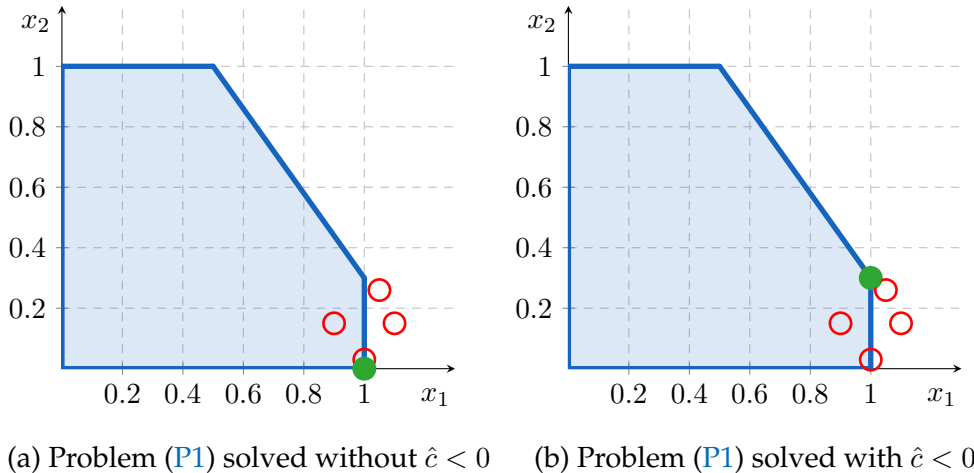


Figure 3: Illustration of example described in Remark 2. Red hollow circles denote the noisy observations, whereas the denoised estimates are depicted by green circles.

Remark 3. So far, we have defined the loss function, i.e. the objective function of (P1), to be the sum of a general norm of the residuals. Typically, a p -norm is used. While the 2-norm seems to work well in most cases, ideally, the norm should be chosen based on the type of noise in the data. Specifically, if the observations are known to be prone to outliers, the 1-norm can be used as it is known to be robust against outliers; whereas, if the data is highly accurate and the hypothesis

of a linear objective function is being tested, one might want to use the ∞ -norm to minimize the worst-case residuals.

4 A Sequential Decomposition Algorithm for (P1)

The phase-1 problem (P1) is an MILP or MINLP (depending on the choice of loss function) whose size increases with the number of experiments, inducing computational challenges when the dataset is large. Therefore, in the following, under the assumption that Condition 1 holds everywhere, we present an exact decomposition algorithm that can substantially reduce the computation time in large instances.

Notice that (P1) has a clear decomposable structure. Specifically, \hat{c} act as linking variables such that with fixed \hat{c} , (P1) decomposes into $|\mathcal{I}|$ independent subproblems, one for each experiment. However, due to the nonconvex nature of the subproblems, traditional Benders-type decomposition methods cannot be directly applied to solve the problem to provable optimality. Instead, we develop a decomposition method that is more akin to Lagrangean decomposition but exploits the structure of the problem such that the exact optimal solution can be obtained. We start by presenting some properties of (P1) that form the basis of our solution algorithm.

Lemma 3. *Given a set of experiments $\{1, \dots, N\}$, let (P1)_{*i*} with $i \leq N$ be an instance of (P1) with $\mathcal{I} = \{i\}$ and let its optimal value be \bar{r}_i^* . If r^* denotes the optimal value of (P1) with $\mathcal{I} = \{1, \dots, N\}$, then*

$$\sum_{i=1}^N \bar{r}_i^* \leq r^*.$$

Proof. Consider the following reformulation of (P1):

$$\begin{aligned}
& \underset{\hat{c}, \hat{x}, s, \lambda, z, w, \hat{c}^+, \hat{c}^-}{\text{minimize}} && \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}_i} \|x_{ij} - \hat{x}_i\| \\
& \text{subject to} && \hat{c}_i = \hat{c}_{i+1} \quad \forall i \in \{1, \dots, N-1\} \\
& && \hat{c}_i + A_i^\top \lambda_i = 0 \quad \forall i \in \mathcal{I} \\
& && A_i \hat{x}_i + s_i = b_i \quad \forall i \in \mathcal{I} \\
& && \lambda_i \leq M z_i \quad \forall i \in \mathcal{I} \\
& && s_i \leq M(e - z_i) \quad \forall i \in \mathcal{I} \\
& && e^\top z_i \geq n \quad \forall i \in \mathcal{I} \\
& && \hat{c}_i = \hat{c}_i^+ - \hat{c}_i^- \quad \forall i \in \mathcal{I} \\
& && \hat{c}_i^+ \leq w_i \quad \forall i \in \mathcal{I} \\
& && \hat{c}_i^- \leq e - w_i \quad \forall i \in \mathcal{I} \\
& && e^\top (\hat{c}_i^+ + \hat{c}_i^-) = 1 \quad \forall i \in \mathcal{I} \\
& && \hat{x}_i \in \mathbb{R}^n, s_i \in \mathbb{R}_+^m, \lambda_i \in \mathbb{R}_+^m, z_i \in \{0, 1\}^m \quad \forall i \in \mathcal{I} \\
& && \hat{c}_i \in \mathbb{R}^n, \hat{c}_i^+ \in \mathbb{R}_+^n, \hat{c}_i^- \in \mathbb{R}_+^n, w_i \in \{0, 1\}^n \quad \forall i \in \mathcal{I},
\end{aligned} \tag{P1'}$$

where we replace the variable \hat{c} in (P1) with N disaggregated variables denoted by \hat{c}_i for each experiment $i \in \mathcal{I}$. We further add the first set of constraints, which equate all \hat{c}_i to ensure (P1')

is equivalent to (P1). Consider a relaxation of (P1') obtained by removing these first set of constraints. With this change, the problem decomposes into N independent subproblems, one for each experiment $i \in \mathcal{I}$, which is precisely (P1) _{i} . It follows that the optimal value of the relaxed problem is $\sum_{i=1}^N \bar{r}_i^*$, and since it is a relaxation of (P1') and therefore also a relaxation of (P1),

$$\sum_{i=1}^N \bar{r}_i^* \leq r^*. \quad \square$$

Corollary 2. Let (P1) _{$[i]$} be an instance of (P1) with $\mathcal{I} = \{1, \dots, i\}$ and $\bar{r}_{[i]}^*$ be its optimal value. Then, $\bar{r}_{[i-1]}^* + \bar{r}_i^* \leq \bar{r}_{[i]}^*$.

Proof. Corollary 2 follows from Lemma 3. □

Our sequential decomposition approach for (P1) is a direct consequence of Corollary 2. We first provide an intuitive description of the method using the illustrative example shown in Figure 4, which involves three experiments. The main idea is to solve (P1) *sequentially*, i.e. one experiment by one experiment, instead of directly solving the full-space problem considering all experiments. With only one experiment, (P1) reduces to the problem of projecting the noisy observations onto the vertex that minimizes the loss. When the level of noise is small, it is likely that this already results in the solution that is optimal for the full problem, as it is the case for experiments 1 and 2 in Figure 4. In experiment 3, however, the level of noise is so high that the projection yields the wrong vertex. A solution to the full problem involving all experiments has to provide a \hat{c} that renders all \hat{x}_i optimal, which implies that the resulting admissible set $\hat{\mathcal{C}}$ must not only contain 0. As shown at the top right of Figure 4, while the true cost vector c lies in the intersection of the two cones resulting from the correct projections, the intersection of all three cones including the one associated with the incorrect projection is $\{0\}$. When $\hat{\mathcal{C}} = \{0\}$, we solve (P1) considering all experiments up to that point jointly in order to correct the projections such that the resulting $\hat{\mathcal{C}}$ is a proper polyhedral cone (see bottom left of Figure 4).

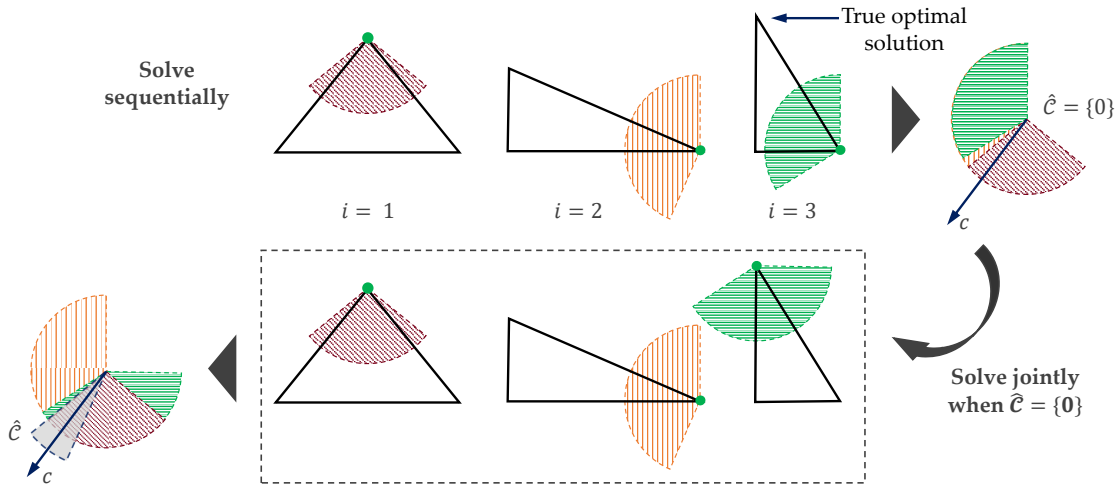


Figure 4: Depiction of the proposed decomposition algorithm. Note that only the vertex projections are shown (filled green circles), the corresponding noisy observations are omitted for the sake of clearer visualization.

Detecting such infeasible projections is crucial for the proposed decomposition algorithm, and it turns out that this can be accomplished by solving a very efficient feasibility problem, which we establish in the following proposition.

Proposition 2. *Let $(\hat{x}_1, \dots, \hat{x}_{\ell-1})$ and \hat{x}_ℓ be feasible solutions to problems $(\mathbf{P1})_{[\ell-1]}$ and $(\mathbf{P1})_\ell$, respectively. Then, $(\hat{x}_1, \dots, \hat{x}_\ell)$ is a feasible solution to $(\mathbf{P1})_{[\ell]}$ if the following problem is feasible:*

$$\begin{aligned}
& \underset{\hat{c}, \hat{c}^+, \hat{c}^-, \gamma, w}{\text{minimize}} && 0 \\
& \text{subject to} && \hat{c} = \sum_{t \in \mathcal{T}(\hat{x}_i)} \gamma_{it} (-a_{it}^\top) \quad \forall i \in \{1, \dots, \ell\} \\
& && \hat{c} = \hat{c}^+ - \hat{c}^- \\
& && \hat{c}^+ \leq w \\
& && \hat{c}^- \leq e - w \\
& && e^\top (\hat{c}^+ + \hat{c}^-) = 1 \\
& && \gamma_{it} \geq 0 \quad \forall i \in \{1, \dots, \ell\}, t \in \mathcal{T}(\hat{x}_i) \\
& && \hat{c} \in \mathbb{R}^n, \hat{c}^+ \in \mathbb{R}_+^n, \hat{c}^- \in \mathbb{R}_+^n, w \in \{0, 1\}^n.
\end{aligned} \tag{FP}_\ell$$

Proof. By construction, a feasible $(\mathbf{FP})_\ell$ implies that $\{\hat{x}_1, \dots, \hat{x}_\ell\}$ are such that there exists a \hat{c} in $\bigcap_{i \in \mathcal{I}} \text{cone}(\{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i)})$ with $\|\hat{c}\|_1 = 1$. This indicates that for the given solutions to $(\mathbf{P1})_{[\ell-1]}$ and $(\mathbf{P1})_\ell$, one can find a common \hat{c} that is feasible in both problems, which in turn implies that $(\hat{x}_1, \dots, \hat{x}_\ell)$ is feasible in $(\mathbf{P1})_{[\ell]}$. \square

Corollary 3. *Let $(\hat{x}_1, \dots, \hat{x}_{\ell-1})$ and \hat{x}_ℓ be optimal solutions to problems $(\mathbf{P1})_{[\ell-1]}$ and $(\mathbf{P1})_\ell$, respectively. Then, $(\hat{x}_1, \dots, \hat{x}_\ell)$ is an optimal solution to $(\mathbf{P1})_{[\ell]}$ if $(\mathbf{FP})_\ell$ is feasible.*

Proof. Due to Corollary 2, the sum of the optimal values of $(\mathbf{P1})_{[\ell-1]}$ and $(\mathbf{P1})_\ell$ provide a lower bound on the optimal value of $(\mathbf{P1})_{[\ell]}$, i.e. $\bar{r}_{[\ell-1]}^* + \bar{r}_\ell^* \leq \bar{r}_{[\ell]}^*$. If $(\mathbf{FP})_\ell$ is feasible, $(\hat{x}_1, \dots, \hat{x}_\ell)$ is feasible in $(\mathbf{P1})_{[\ell]}$ due to Proposition 2 and therefore yields an upper bound, i.e. $\bar{r}_{[\ell-1]}^* + \bar{r}_\ell^* \geq \bar{r}_{[\ell]}^*$. As the lower and upper bounds coincide, it follows that $(\hat{x}_1, \dots, \hat{x}_\ell)$ is an optimal solution to $(\mathbf{P1})_{[\ell]}$. \square

Although $(\mathbf{FP})_\ell$ is an MILP, its number of binary variables does not change with the number of data points, which allows it to remain tractable for instances with many experiments. Therefore, in the algorithm, after every experiment ℓ , we solve $(\mathbf{FP})_\ell$ to check the validity of the partial solution obtained up to that point for problem $(\mathbf{P1})_{[\ell]}$. If at some point, $(\mathbf{FP})_\ell$ is found infeasible, we solve the full problem $(\mathbf{P1})_{[\ell]}$. Here, it is important to note that since the solution obtained for $(\mathbf{P1})_{[\ell-1]}$ in the previous iteration is guaranteed to be feasible for $(\mathbf{P1})_{[\ell]}$, it can be used to warm-start the solver, which is another crucial factor with regard to the computational efficiency of the algorithm. In practice, with low level of noise or a sufficiently large number of samples $|\mathcal{J}_i|$ for each experiment $i \in \mathcal{I}$, one can expect the single-experiment projection to be accurate most of the time. As a result, one may only need to solve a small number of problems involving multiple experiments. The pseudocode for the complete algorithm for solving (\mathbf{IOP}) with this sequential decomposition approach is shown in Algorithm 2.

Algorithm 2 Decomposition-based algorithm for solving (IOP)

```
1: initialize:  $\hat{\mathcal{C}} \leftarrow \mathbb{R}^n$ 
2: for all  $\ell \in \{1, \dots, N\}$  do
3:   solve (P1) $_{\ell}$ , obtain  $\hat{x}_{\ell}^*$ 
4:    $\hat{\mathcal{C}} \leftarrow \hat{\mathcal{C}} \cap \text{cone} \left( \{-a_{\ell t}\}_{t \in \mathcal{T}(\hat{x}_{\ell}^*)} \right)$ 
5:   solve (FP) $_{\ell}$ 
6:   if (FP) $_{\ell}$  is infeasible then
7:     use Algorithm 1 to solve (P1) $_{[\ell]}$ , warm-start with  $(\hat{x}_1^*, \dots, \hat{x}_{\ell-1}^*)$ 
8:      $\hat{\mathcal{C}} \leftarrow \bigcap_{i \in \mathcal{I}} \text{cone} \left( \{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i^*)} \right)$ 
9:   end if
10: end for
11: solve (P2), obtain  $\hat{c}^*$ 
```

Online Inverse Optimization So far, we have assumed that all data are available prior to the inverse optimization process. However, in many situations, it is reasonable to expect the data from different experiments to become available at different time points. This has motivated the development of online learning frameworks for inverse optimization where the cost estimate is updated as new data arrive (Bärmann et al., 2017; Dong et al., 2018; Shahmoradi & Lee, 2019). A naive extension of our two-phase framework for online learning can be to repeatedly solve (IOP) with the new data added to it. However, note that Algorithm 2 provides a much more efficient mechanism to address this problem, and can be adapted as a method for inverse optimization in an online setting. We include such an adaptation in Appendix A-1.

5 Online Adaptive Sampling

In this section, we introduce a novel adaptive sampling approach that can be used to reduce the number of experiments required to obtain a good cost estimate in an online setting. Our method derives from the geometrical understanding of (IOP) discussed in the previous sections. As a motivating example, let us revisit the case with three experiments shown in Figure 4. Notice that the admissible set $\hat{\mathcal{C}}$ obtained with the first two experiments is the same as the one obtained with all three experiments. The reason is that the cone associated with the projected solution from experiment 3 is a superset of the cone from experiment 1; hence, the intersection of cones does not change with the third experiment. This means that experiment 3 cannot help improve the cost estimate or, more precisely, our confidence in the cost estimate, which only increases if the size of $\hat{\mathcal{C}}$ decreases. Therefore, the goal here is to use the current best cost estimate to design subsequent experiments in a way such that the size of $\hat{\mathcal{C}}$ will likely be further reduced. In the following, we formulate the adaptive sampling problem as an optimization problem, and devise an efficient heuristic method to solve it.

5.1 Mathematical Formulation

Consider the point in an online inverse optimization process at which $\ell - 1$ experiments are completed and the problem is to choose input parameters $(A_\ell, b_\ell) \in \Pi$ for the ℓ th experiment. Let \mathcal{I} be the set of the first $\ell - 1$ experiments, and let $\widehat{\mathcal{C}}$ be the admissible set obtained using these experiments. Suppose $\widehat{\mathcal{C}}_\ell$ denotes the (unknown) cone formed by the active constraints of the FOP with its feasible region represented by the polyhedron $\{\tilde{x} : A_\ell \tilde{x} \leq b_\ell\}$. The goal as illustrated in the above example is to choose A_ℓ and b_ℓ such that $\widehat{\mathcal{C}} \cap \widehat{\mathcal{C}}_\ell \subset \widehat{\mathcal{C}}$. While it is impossible to determine $\widehat{\mathcal{C}}_\ell$ without knowing c exactly, we make use of a randomly sampled vector $\tilde{c}_{\ell-1}$ from $\widehat{\mathcal{C}}$ to predict which constraints may become active for given sets of potential input parameters.

Assuming we have an accurate estimate of $\widehat{\mathcal{C}}_\ell$, we utilize the fact that $\widehat{\mathcal{C}} \cap \widehat{\mathcal{C}}_\ell \subset \widehat{\mathcal{C}}$ if and only if we can find a point that belongs to $\widehat{\mathcal{C}}$ but not to $\widehat{\mathcal{C}}_\ell$. As we illustrate in Figure 5, the existence of such a point can be determined by obtaining the minimum-distance projections of all points in $\widehat{\mathcal{C}}$ on $\widehat{\mathcal{C}}_\ell$ and searching for one that yields a nonzero projection distance. In the figure, \bar{y} denotes the minimum-distance projection of a point y from $\widehat{\mathcal{C}}$ on $\widehat{\mathcal{C}}_\ell$. In case of Figure 5a, all points $y \in \widehat{\mathcal{C}}$ are such that $\bar{y} = y$ and hence, the intersection of the two cones does not result in a reduction in the size of the admissible set. In Figure 5b, however, one can find a y and its corresponding \bar{y} with $\|y - \bar{y}\| > 0$ and therefore, the intersection of the two cones is a strict subset of $\widehat{\mathcal{C}}$.

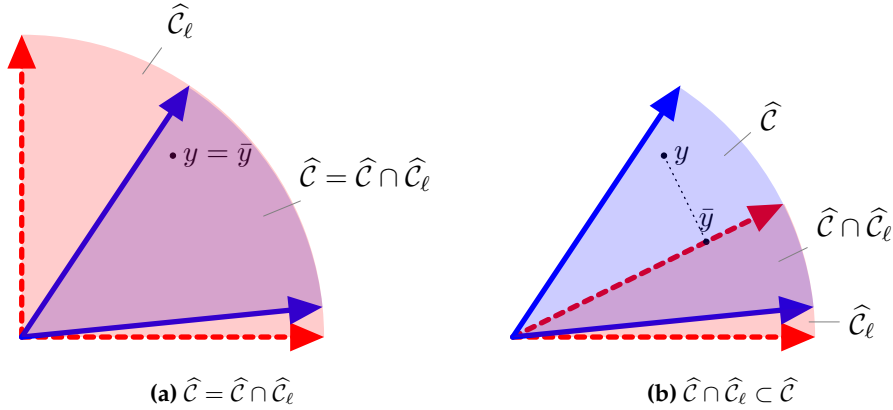


Figure 5: The cone formed by the blue vectors represents $\widehat{\mathcal{C}}$ whereas the red dashed vectors form the cone $\widehat{\mathcal{C}}_\ell$. The dotted line in the second case shows the minimum 2-norm projection of y on $\widehat{\mathcal{C}}_\ell$.

Given $\widehat{\mathcal{C}}$ and $\widehat{\mathcal{C}}_\ell$, we find the $y \in \widehat{\mathcal{C}}$ with the maximum minimum-projection distance to $\widehat{\mathcal{C}}_\ell$ and use this distance as a measure for the difference between $\widehat{\mathcal{C}}$ and $\widehat{\mathcal{C}} \cap \widehat{\mathcal{C}}_\ell$. Hence, as we wish to reduce the size of the admissible set as much as possible, we need to choose A_ℓ and b_ℓ for which this distance is the largest among all $(A_\ell, b_\ell) \in \Pi$. We formulate the problem of finding such A_ℓ and b_ℓ as the following optimization problem:

$$\begin{aligned} & \underset{y, \gamma, z, A_\ell, b_\ell, \tilde{x}_\ell}{\text{maximize}} \quad \min_{\eta, \bar{y}, \tilde{\gamma}} \left\{ \eta : \eta = \|y - \bar{y}\|_1, \bar{y} = \sum_{k \in \mathcal{K}} \tilde{\gamma}_k z_k a_{\ell k}, \tilde{\gamma}_k \geq 0 \forall k \in \mathcal{K} \right\} & \text{(ASPa)} \end{aligned}$$

$$\begin{aligned} & \text{subject to} \quad y = \sum_{t \in \mathcal{T}(\tilde{x}_i)} \gamma_{it} \left(-\frac{a_{it}}{\|a_{it}\|_2} \right) \quad \forall i \in \mathcal{I} & \text{(ASPb)} \end{aligned}$$

$$0 \leq \gamma_{it} \leq \epsilon \quad \forall i \in \mathcal{I}, t \in \mathcal{T}(\hat{x}_i) \quad (\text{ASPC})$$

$$(A_\ell, b_\ell) \in \Pi \quad (\text{ASPD})$$

$$z_k = \begin{cases} 1, & \text{if } k \in \mathcal{T}(\hat{x}_\ell) \text{ where } \hat{x}_\ell \in \arg \min_{\tilde{x} \in \mathbb{R}^n} \{ \tilde{c}_{\ell-1}^\top \tilde{x} : A_\ell \tilde{x} \leq b_\ell \} \\ 0, & \text{otherwise} \end{cases} \quad \forall k \in \mathcal{K}, \quad (\text{ASPE})$$

where \mathcal{K} is the set of constraints for $A_\ell x \leq b_\ell$. Constraints (ASPB) state that y is a point in $\hat{\mathcal{C}}$. Note that the vectors forming cone $\hat{\mathcal{C}}_i$ for each $i \in \mathcal{I}$ have been normalized to ensure that the lengths of the vectors do not bias the value of η . The upper bound ϵ in (ASPC) ensures that the problem remains bounded. The value of ϵ can be chosen arbitrarily and does not affect the choice of optimal A_ℓ and b_ℓ for the problem. Constraint (ASPD) imposes the restriction that input parameters A_ℓ and b_ℓ must be chosen from the set Π . As stated in constraints (ASPE), z_k is a binary variable that equals 1 if $\tilde{c}_{\ell-1}$ predicts that the k th constraint will be active. Finally, the objective of (ASP) is to maximize the minimum 1-norm distance between y and \bar{y} , with \bar{y} being constrained to lie in the cone $\hat{\mathcal{C}}_\ell$. We denote the distance between y and \bar{y} by η .

5.2 A Heuristic Solution Algorithm for (ASP)

Problem (ASP) is a bilevel optimization problem with lower-level problems embedded in the objective function (ASPA) and constraints (ASPE). Since these embedded problems are LPs, (ASP) can be reformulated into a single-level problem, similar to the bilevel problems considered in previous sections. We present two such single-level reformulations in Appendix A-2. Both of them take the forms of nonconvex MINLPs that are only tractable for relatively small problems. Next, we outline a heuristic approach to solve this problem.

Instead of searching through the entire set Π for the best set of input parameters, the proposed algorithm, Algorithm 3, simplifies the problem by limiting the search space to S randomly sampled sets of input parameters. It then solves S instances of (ASP), each corresponding to one of the parameter sets to choose the one that results in the largest η . Here, since the number of possible candidates for A_ℓ and b_ℓ is finite with their values being explicitly known, binary vector z for each of them can be calculated before (ASP) is solved. This eliminates the lower-level problems in (ASPE). Furthermore, in this algorithm, since the S instances of (ASP) are independent of each other, they can be solved in parallel. This enables the use of a large S to find a heuristic solution as close to the global optimum of (ASP) as possible while still keeping the problem tractable.

Algorithm 3 A heuristic algorithm for solving (ASP)

- 1: **for** $i \in \{1, \dots, S\}$ **do**
 - 2: sample $(A_i, b_i) \in \Pi$
 - 3: compute z_i such that $z_{ik} = \begin{cases} 1, & \text{if } k \in \mathcal{T}(\hat{x}_i) \text{ where } \hat{x}_i \in \arg \min_{\tilde{x} \in \mathbb{R}^n} \{ \tilde{c}_{i-1}^\top \tilde{x} : A_i \tilde{x} \leq b_i \} \\ 0, & \text{otherwise} \end{cases}$
 - 4: solve (ASP) with $A_\ell = A_i$, $b_\ell = b_i$, and $z = z_i$, obtain η_i
 - 5: **end for**
 - 6: **return** A_i, b_i for which $\eta_i = \max_{i' \in \{1, \dots, S\}} \eta_{i'}$
-

5.3 Effect of Restrictions on the Design of Experiments

Problem (ASP) is valid regardless of the restrictions on the input parameters A and b , i.e. the set Π . However, it is still worth highlighting the difference between being able to vary the constraint matrix A versus being restricted to only varying the right-hand-side (RHS) parameters b . While the first case is more general (we include two case studies in the next section), the second scenario is also not uncommon. In several transportation network flow problems, the values in A are restricted to 0 and 1. These values determine the structure of the network and cannot be changed. The RHS parameters b , however, usually refer to quantities such as capacities of individual arcs or demands at the nodes and can often be controlled. Notice that the primary requirement for adaptive or even random sampling to reduce the size of $\hat{\mathcal{C}}$ is to find cones in successive experiments with extreme rays pointing in directions different from the ones already found. With A being allowed to vary, this can be easily achieved by appropriately varying the parameters of the active constraints under the given restrictions Π . However, when only b can vary, the only hope to reduce the size of the admissible set is in finding a set of input parameters for which a different set of constraints become active. We further illustrate this point through the following example.

Example 1. Let the (FOP) for the first experiment be $\max_{x_1, x_2} \{x_1 + 2x_2 : 2x_1 + 3x_2 \leq 3, x_1 \leq 1, x_2 \leq 1, x_1 \geq 0, x_2 \geq 0\}$, for which the feasible region is depicted in the first subfigure of Figure 6. Suppose the constraints are indexed by $\mathcal{K} = \{1, 2, 3, 4, 5\}$ and the set Π allows $2.5 \leq b_1 \leq 4$, keeping all other parameters fixed. For simplicity, we assume that the noisy data obtained is such that the correct vertex can always be recovered by solving (P1) for each of the experiments individually. Therefore, $\hat{\mathcal{C}}_1 = \text{cone}(\{(-1, 0), (2, 3)\})$. Now, for the next experiment to result in a cone with $\eta > 0$, we have to choose $b_1 > 3$, which will yield a different set of active constraints. For any other choice of b_1 , $\hat{\mathcal{C}}_2$ will be the same as $\hat{\mathcal{C}}_1$. In Figure 6, we show the case when b_1 for the second experiment is chosen to be 4. This b_1 results in a $\hat{\mathcal{C}}_2 = \text{cone}(\{(0, 1), (2, 3)\})$, thereby effectively reducing the size of the admissible set and increasing the confidence in the point estimate selected from this set.

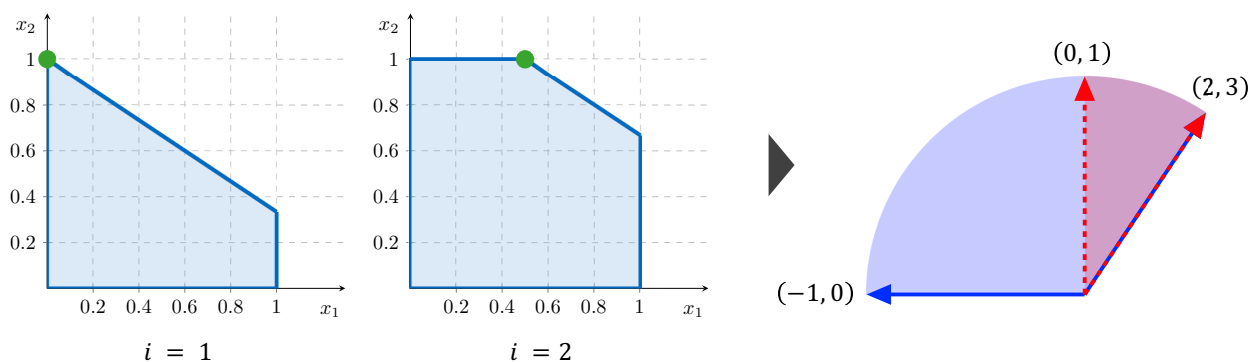


Figure 6: Illustration of the problem described in Example 1. The vectors are normalized to the same length. Filled green circles depict the denoised estimates of the noisy observations (omitted).

6 Computational Case Studies

In this section, we apply the proposed data-driven inverse linear optimization framework to two case studies, one addressing customer preference learning and the second related to cost estimation for multiperiod production planning. Using synthetic data, we compare the computational performances of Algorithms 1 and 2 and evaluate the impact of adaptive sampling. All model instances were implemented in Julia v1.3.0 (Bezanson et al., 2017) using the modeling language JuMP v0.18.6 (Dunning et al., 2017). All instances of (P1) were solved with a 1-norm-based objective function, and Gurobi v9.0.2 (Gurobi Optimization, 2020) was applied to solve the resulting MILPs. Nonconvex MINLPs were solved using BARON v19.12.7 (Sahinidis, 1996).

6.1 Customer Preference Learning

We consider the problem of learning customers' preferences given their purchasing decisions. The FOP here is based on the premise that with a limited budget, customers will buy products that maximize their utility. Given the price w_p of each product p and a budget b , the customer is assumed to solve the following LP:

$$\begin{aligned} & \underset{x \in \mathbb{R}_+^n}{\text{maximize}} && \sum_{p \in \mathcal{P}} u_p x_p \\ & \text{subject to} && \sum_{p \in \mathcal{P}} w_p x_p \leq b \\ & && x_p \leq 1 \quad \forall p \in \mathcal{P}, \end{aligned} \tag{12}$$

where $\mathcal{P} = \{1, \dots, n\}$ denotes the set of n products available on the market. The goal is to estimate the unknown utility function coefficients u by observing the changes in the customer's decisions x in response to price fluctuations. Different variations of this IOP have previously been considered by Bärmann et al. (2017) and Dong et al. (2018). Bärmann et al. (2017) consider learning the utility function with deterministic data. Dong et al. (2018) account for noise in the data but pose (12) with a strongly concave utility function. Our design of this case study closely follows the scheme presented by Bärmann et al. (2017), but uses noisy data to learn the utility function coefficients.

For each instance of the IOP, the training data are generated as follows. We first create an arbitrary utility vector $u \in \mathbb{R}^n$ by sampling its individual elements from the uniform distribution $\mathcal{U}(1, 1000)$, and normalize it to make its 1-norm equal to 1. We then sample a set of price vectors w_i , which are the input parameters for each $i \in \mathcal{I}$ such that $w_{ip} \sim \mathcal{U}(50, 150)$ for every $p \in \mathcal{P}$. The budget b is set to $0.6 \sum_{p \in \mathcal{P}} w_{1p}$ for all experiments. Next, keeping the utility vector the same, we solve these $|\mathcal{I}|$ instances of (12) to obtain the optimal decisions x_i^* . We then generate the noisy datasets \mathcal{J}_i for each $i \in \mathcal{I}$ by distorting the true optimal solution such that $x_{ij} = x_i^* + \gamma$, where $\gamma \sim \mathcal{N}(0, \sigma^2 \mathbb{I})$.

In this study, we consider FOPs of varying dimensionality n but limit the number of experiments to 100 in all cases. We also consider datasets with varying levels of noise by changing the value of σ . Once n and σ are fixed, the size of the sets \mathcal{J}_i is kept the same for all $i \in \mathcal{I}$, i.e. $|\mathcal{J}_i| = J$ for all $i \in \mathcal{I}$. A specific case is hence represented by n , σ , and J , and we solve ten random instances of each case (generated using the scheme described above). Following Remark 2, to

increase the robustness of our two-phase algorithm against such a large number of polyhedral geometries and different levels of noise, we introduce an additional set of constraints in the phase-1 problem enforcing $\hat{u}_p \geq 10^{-6}$ for all $p \in \mathcal{P}$.

6.1.1 Computational Performance

The results comparing the computational performance of Algorithms 1 and 2 are summarized in Table 1. All instances were solved with a time limit of 7,200 s utilizing 24 cores on the Mesabi cluster of the Minnesota Supercomputing Institute (MSI). For each of the algorithms, the table lists the number of instances (out of ten) that were solved to optimality. In all instances that could not be solved with Algorithm 1, we find that the solver could not find even a single feasible solution within the given time; hence the optimality gaps are not reported. For the decomposition-based algorithm, an unsolved instance is one where not all the 100 experiments could be processed in the given time but it still yields an estimate for the cost coefficients. Nonetheless, for the comparison of the two solution methods, we consider a run resulting in a partial solution as a failed run. Table 1 shows the median, maximum, and minimum computation times for both algorithms. For Algorithm 1, it is the time required to solve an instance without implementing the integer cuts in Phase 1. We find that adding the integer cuts to identify multiple optimal solutions was only required in 2 of the 180 ($\sim 1\%$) instances considered in this study, confirming our assertion that a violation of Condition 1 is rare. Lastly, note that Algorithm 2 has an additional column labeled “median # of re-solves” that shows the median value of the number of times the feasibility problem $(\text{FP})_\ell$ was found infeasible, which then required the problem involving all experiments up to that point to be re-solved.

From the data in Table 1, one can observe that the two algorithms solve the IOP in comparable times when the level of noise is low. However, as the size of the problem or the level of noise increases, Algorithm 2 starts outperforming Algorithm 1. The difference in their performance is especially apparent from the numbers of instances solved and the maximum computation times. For example, in the arguably most difficult case with $n = 100$, $\sigma = 0.1$, and $J = 20$, Algorithm 1 was not able to solve any of the given instances, while Algorithm 2 solved eight out of ten. Irrespective of the solution algorithm, increasing J has a seemingly counterintuitive effect of making the problem easier to solve. This is because increasing the number of samples for an experiment merely increases the number of terms in the objective function of (P1) while making it easier for the problem to find the correct vertex to project onto. For Algorithm 2, one can also see that increasing J reduces the likelihood of incorrect projections when processing the data for individual experiments separately. This is especially helpful in situations where computing infrastructure is a limitation as large amount of data has a significantly higher memory requirement. Overall, the results indicate that Algorithm 2 dominates over Algorithm 1 when it comes to solving more difficult instances of (IOP) with data of high dimensionality and level of noise.

6.1.2 Prediction Error

Phase 2 of our two-phase algorithm requires a reference \bar{u} to yield an estimate \hat{u} . Ideally, this reference is based on some prior intuition about the unknown utility function, but here we obtain \hat{u} using a randomly generated \bar{u} . The goal is to test the capability of this estimate in generating

n	σ	J	Algorithm 1				Algorithm 2				median # of re-solves
			# of instances solved	computation time (s)			# of instances solved	computation time (s)			
				median	max	min		median	max	min	
25	0.01	5	10	33	155	6	10	277	317	253	2
		250	10	73	90	67	10	266	271	264	0
	0.05	10	10	176	264	54	10	283	789	262	6
		250	10	88	965	74	10	268	347	264	0
	0.1	20	10	306	680	22	10	315	499	273	9.5
		250	10	180	5,563	79	10	298	687	265	1
50	0.01	5	10	88	931	7	10	484	638	331	1.5
		250	10	220	1,724	200	10	423	717	414	0
	0.05	10	10	1,037	5,874	430	10	673	6,897	409	7.5
		250	7	314	439	215	10	626	1,617	417	1.5
	0.1	20	8	2,707	7,081	715	10	817	1,841	504	11.5
		250	8	875	6,035	265	10	705	1,651	510	2
100	0.01	5	9	407	488	29	10	428	1,042	362	2
		250	10	489	2,267	419	10	479	1,648	451	0
	0.05	10	3	2,707	2,851	1,902	9	1,021	2,032	586	9
		250	9	2,232	6,610	526	10	865	5,537	455	1.5
	0.1	20	0	n/a	n/a	n/a	8	1,656	2,933	935	12
		250	4	1,089	2,561	542	8	1,337	6,890	565	4.5

Table 1: Comparison of computational performances of Algorithms 1 and 2 on an IOP based on random instances of (12). Reported computation times only consider the instances that were solved to optimality.

reliable predictions on unseen datasets in the worst-case scenario where no prior information about the missing parameters is available. To perform this assessment, along with every instance of training data, we also generate a test dataset of 100 experiments. This test data consists of (w, x^*) pairs, where w -values are generated in the same manner as for the training data, and x^* are the corresponding true optimal solutions obtained using the same u as the one used to generate the noisy training data. Once a \hat{u} has been found, we use it to solve the problems in the test dataset and evaluate the prediction error as the fraction of incorrect predictions.

We show the prediction error results for a few selected cases in Figure 7. Here, instead of using just the final estimate obtained with all the $|\mathcal{I}|$ experiments, we make use of the online algorithm (Algorithm 4 in Appendix A-1) to show how the prediction error evolves with the addition of new experiments to the training set. As expected, the prediction error generally decreases with the number of experiments, and problems of higher dimensionality require more experiments to reach the same prediction accuracy. Also, the figures show an additional benefit of a large J apart from making (P1) easier to solve. In Figures 7a and 7b, where the size of J is kept small, the curves are “spiky” showing a local increase in prediction error. Recall that we solve $(\text{FP})_\ell$ after

every experiment ℓ to confirm if the the solution obtained with the first ℓ experiments ($\ell \leq |\mathcal{I}|$) still holds for $(P1)_{[\ell]}$. Therefore, these local spikes are a consequence of insufficient sampling which results in the violation of our primary assumption that the vertex with minimum loss on the given data is the “correct” vertex. As seen in Figure 7c, once an adequate number of samples are used, the likelihood of observing these peaks reduces significantly. We expect that in situations where the quality of data is uncertain and sampling is limited, adding a pre-processing step to remove large outliers from the dataset can be a way to prevent the model from making wrong estimates.

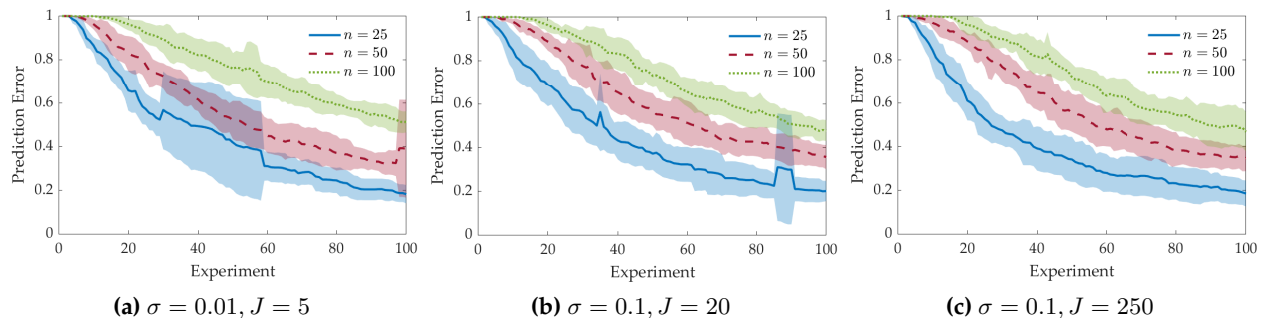


Figure 7: Change in prediction error as experiments are added. Lines show mean values across all solved instances, shaded areas indicate one standard deviation around the mean. For all the cases, prediction error data for individual instances are available in Appendix A-3.

6.1.3 Adaptive Sampling

We also apply our adaptive sampling strategy to this example. Following the data generation scheme used for random sampling, we define the set Π as allowing any w for which $50 \leq w_p \leq 150$ for any $p \in \mathcal{P}$ while keeping the other constraint parameters fixed. Irrespective of the size of (IOP), we find that BARON struggles to find even a feasible solution for (ASP) when we try to solve it exactly. Therefore, we solve it here using the heuristic solution algorithm, Algorithm 3. All the S subproblems were solved with a time limit of 100 s. Although BARON is unable to solve the subproblems to optimality, it still finds feasible solutions that show good potential in reducing the size of the admissible set. The significant impact of our adaptive sampling strategy is shown in Figure 8. In all three cases, the use of adaptive sampling results in a more than 50% reduction in the number of experiments required to achieve the same prediction accuracy as obtained from the standard approach using random sampling after 100 experiments. Furthermore, Figure 8a shows the effect of the parameter S in the heuristic solution approach for (ASP). One can see that even a small S results in a large increase in the rate of reduction of the prediction error and moreover, the rate stabilizes very quickly for a rather small S . By performing a similar study on higher-dimensional problems, we find that using an S equal to the dimensionality of the problem n to be a good heuristic to achieve the desired effect with adaptive sampling.

6.2 Cost Estimation for Production Planning

In our second case study, we consider the problem of production planning for a large manufacturing site consisting of multiple processes within an interconnected process network. It is commonly

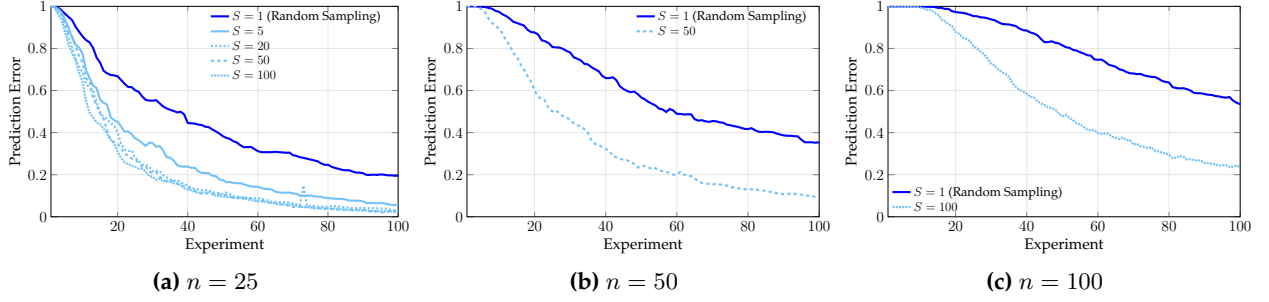


Figure 8: Effect of adaptive sampling on the evolution of prediction error. Lines show mean values across ten instances. Training data for all cases was generated using $\sigma = 0.01$, $J = 30$.

formulated as an LP:

$$\underset{x,y,w}{\text{minimize}} \quad \sum_{h \in \mathcal{H}} \sum_{m \in \mathcal{M}} \left(\sum_{p \in \mathcal{P}} c_{pmh} x_{pmh} + f_{mh} w_{mh} \right) \quad (13a)$$

$$\text{subject to} \quad q_m^{\min} \leq q_m^0 + \sum_{h'=1}^h \left(\sum_{p \in \mathcal{P}} x_{pmh'} + w_{mh'} - d_{mh'} \right) \leq q_m^{\max} \quad \forall m \in \mathcal{M}, h \in \mathcal{H} \quad (13b)$$

$$0 \leq w_{mh} \leq w_{mh}^{\max} \quad \forall m \in \mathcal{M}, h \in \mathcal{H} \quad (13c)$$

$$x_{pmh} = \mu_{pm} y_{ph} \quad \forall p \in \mathcal{P}, m \in \mathcal{M}, h \in \mathcal{H} \quad (13d)$$

$$0 \leq x_{pmh} \leq x_{pmh}^{\max} \quad \forall p \in \mathcal{P}, m \in \mathcal{M}, h \in \mathcal{H}, \quad (13e)$$

where \mathcal{P} , \mathcal{M} , and $\mathcal{H} = \{1, \dots, H\}$ are the sets of processes, materials, and time periods, respectively. The amount of material m produced or consumed (depending on the sign) by process p in time period h is denoted by x_{pmh} . Product demand and additional purchase of a material are denoted by d_{mh} and w_{mh} , respectively. Inventory constraints and restrictions on the amounts purchased are stated in constraints (13b) and (13c), respectively. The structure of the process network is defined by constraints (13d) where y_{ph} denotes the amount of a reference material for process p produced in time period h , and μ_{pm} is a conversion factor that specifies how much of a material m is produced or consumed for one unit of the reference material. According to (13a), the objective is to minimize the total production and purchasing cost while satisfying given product demand. While purchasing prices f_{mh} are readily known, production costs c_{pmh} are often difficult to estimate, which is one major reason why in practice, production planning is still mostly performed manually by human planners (Troutt et al., 2006). Experienced planners have an excellent intuition for the relative differences in costs, but this information is not explicitly expressed in numbers. The goal is to use past production plans, which reflect the planners' decisions, to infer these costs.

Here we conduct a case study using the process network for a petrochemical site with 38 chemicals and 28 processes from Sahinidis et al. (1989) and Zhang et al. (2016). The parameters required to model this network are given in the Appendix A-4. We test our methodology on this problem by generating synthetic training data simulating expert planners' decision making under different operating and market conditions. Here, each data point consists of inputs (μ, d) and the

corresponding decisions (x, y, w) . For different scenarios, the conversion factors μ_{pm} have been assumed to vary between 75% to 100% of their nominal values on account of changing efficiencies of individual processes. Also, the product demand d_{mh} is considered to fluctuate $\pm 10\%$ from its nominal value. For each instance, we generate 100 such scenarios, i.e. $|\mathcal{I}| = 100$. For each of these scenarios, we then generate the respective decisions by solving (13) with the same arbitrarily generated c_{pmh} and f_{mh} values. Since human decision making is often inconsistent, we distort the true optimal solution (x^*, y^*, w^*) as $(x^* + \gamma_1, y^* + \gamma_2, w^* + \gamma_3)$ where $\gamma_i \sim \mathcal{N}(0, \sigma^2)$ for all $i \in \{1, 2, 3\}$.

In this case study, we consider training data instances of three different sizes by varying the number of time periods H . All training data is generated using $\sigma = 3$ and $J = 30$. Each dataset is used to solve (IOP) with both Algorithm 1 and the decomposition-based Algorithm 2. The problem instances were solved with a time limit of 14,400 s using 24 cores on the Mesabi cluster of the MSI. We find that while Algorithm 1 is unable to find even a feasible solution for any of the three cases, Algorithm 2 can solve these problems in less than 10 minutes.

Unlike the previous case study where we assessed the quality of only a point estimate obtained from Phase 2, here we instead focus on the quality of the set $\hat{\mathcal{C}}$. To do this, we sample ten point estimates by solving (P2) with ten random reference cost vectors. We again consider a test dataset of 100 data points consisting of arbitrary (μ, d) and (x^*, y^*, w^*) pairs. We relax the prediction error criteria to obtain a more realistic metric that measures the closeness of the generated predictions to the true optimal solutions; the metric is defined as $d_{\mathcal{V}}(x^*, \hat{x}) = \sum_{v \in \mathcal{V}} \|x_v^* - \hat{x}_v\|_{\infty}$, where \mathcal{V} is the set of test data points. Figure 9a shows the change in the normalized distance metric with the addition of new experiments. One can verify that the spread around the mean value after 100 experiments is fairly small, implying high confidence in the final point estimate irrespective of the quality of the reference. However, this does not discount the importance of a good reference as the admissible set also likely contains the true c which, if used as a reference, would result in a zero $d_{\mathcal{V}}$ even with a single experiment.

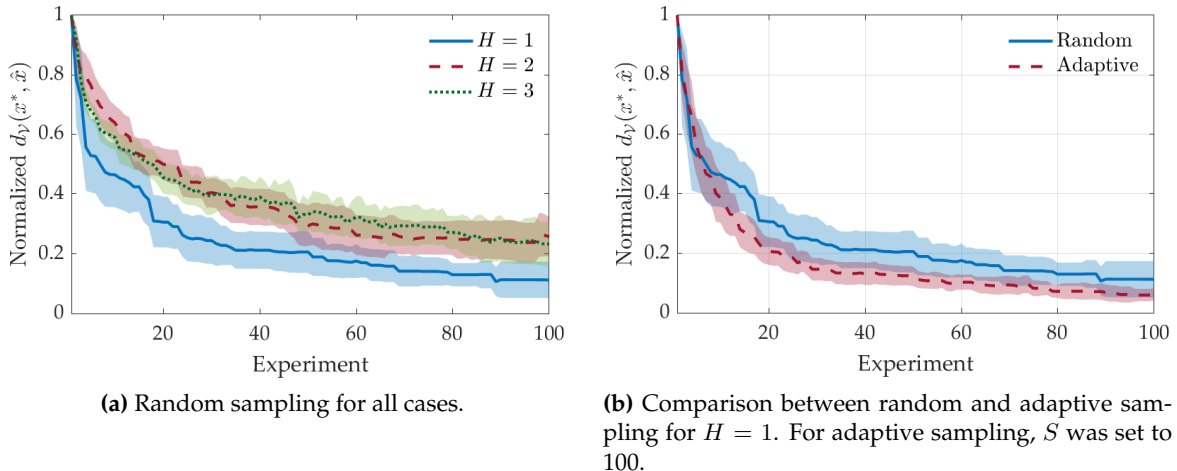


Figure 9: Normalized $d_{\mathcal{V}}(x^*, \hat{x})$ as experiments are added to the training set. Lines show mean values across all instances, shaded areas indicate one standard deviation around the mean. Training data was generated for $\sigma = 3, J = 30$.

Notice that the curve representing the mean value decreases in discrete steps with several flat regions in between two steps. This is a consequence of not all inputs resulting in a reduction in the size of the admissible set. Therefore, we also evaluate the impact of adaptive sampling on mitigating this issue. Due to the larger size of the FOP here compared to (12), we solve the sub-problems in Algorithm 3 with an increased time limit of 200 s. However, even with the increased time limit, the solver struggles to find good feasible solutions for the larger instances with $H = 2$ (132 variables) and $H = 3$ (198 variables). Here, we show the effect of adaptive sampling with the case of $H = 1$ in Figure 9b. As can be observed, adaptive sampling sustains a high rate of decrease in the distance metric for a longer duration compared to naive random sampling. This results in just ~ 52 experiments being required to achieve the same effect as 100 random inputs. Moreover, one can see that the variance around the mean in the case of adaptive sampling is noticeably lower which shows that in addition to requiring less experiments, it also finds estimates with higher confidence levels.

7 Conclusions

In this work, we addressed data-driven inverse linear optimization with noisy observations, for which we introduced a new problem formulation that offers two practical advantages over other existing methods: (i) It allows the recovery of a less restrictive and generally more appropriate admissible set of cost estimates by assuming that the optimal solutions of the FOP lie at the vertices of the feasible region. (ii) Instead of randomly choosing a point estimate from the admissible set, it makes use of a reference cost vector to choose the cost estimate that most resembles the user’s prior belief.

An exact two-phase algorithm was developed to solve the IOP, and we further proposed an efficient decomposition algorithm and an adaptive sampling method that are especially suited for an online inverse optimization setting. Results from extensive computational experiments based on two case studies show that the proposed methods are effective in significantly reducing both the computation time and data requirement for generating cost estimates with a reasonably low prediction error on unseen datasets.

Acknowledgments

The authors gratefully acknowledge the Minnesota Supercomputing Institute (MSI) at the University of Minnesota for providing resources that contributed to the research results reported in this paper. RG acknowledges financial support from a departmental fellowship sponsored by 3M.

References

- Ahuja, R. K. & Orlin, J. B. (2001). Inverse optimization. *Operations Research*, 49(5), 771–783.
- Arechavaleta, G., Laumond, J. P., Hicheur, H., & Berthoz, A. (2008). An optimality principle governing human walking. *IEEE Transactions on Robotics*, 24(1), 5–14.

- Aswani, A., Shen, Z. J. M., & Siddiq, A. (2018). Inverse optimization with noisy data. *Operations Research*, 66(3), 870–892.
- Babier, A., Chan, T. C., Lee, T., Mahmood, R., & Terekhov, D. (2018). An ensemble learning framework for model fitting and evaluation in inverse linear optimization. *arXiv preprint arXiv:1804.04576*.
- Bärman, A., Pokutta, S., & Schneider, O. (2017). Emulating the expert: Inverse optimization through online learning. *34th International Conference on Machine Learning, ICML 2017, 1*, 632–646.
- Beil, D. R. & Wein, L. M. (2003). An inverse-optimization-based auction mechanism to support a multiattribute RFQ process. *Management Science*, 49(11), 1529–1545.
- Bertsimas, D., Gupta, V., & Paschalidis, I. C. (2012). Inverse optimization: A new perspective on the Black-Litterman model. *Operations Research*, 60(6), 1389–1403.
- Bertsimas, D., Gupta, V., & Paschalidis, I. C. (2015). Data-driven estimation in equilibrium using inverse optimization. *Mathematical Programming*, 153(2), 595–633.
- Bezanson, J., Edelman, A., Karpinski, S., & Shah, V. B. (2017). Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1), 65–98.
- Birge, J. R., Hortaçsu, A., & Pavlin, J. M. (2017). Inverse optimization for the recovery of market structure from market outcomes: An application to the MISO electricity market. *Operations Research*, 65(4), 837–855.
- Bulut, A. & Ralphs, T. K. (2016). On the complexity of inverse mixed integer linear optimization. Technical report, Technical Report 15T-001-R3; Lehigh University: Bethlehem, PA, USA.
- Burgard, A. P. & Maranas, C. D. (2003). Optimization-based framework for inferring and testing hypothesized metabolic objective functions. *Biotechnology and Bioengineering*, 82(6), 670–677.
- Burton, D. & Toint, P. L. (1992). On an instance of the inverse shortest paths problem. *Mathematical Programming*, 53(1-3), 45–61.
- Chan, T. C., Craig, T., Lee, T., & Sharpe, M. B. (2014). Generalized inverse multiobjective optimization with application to cancer therapy. *Operations Research*, 62(3), 680–695.
- Chan, T. C. & Lee, T. (2018). Trade-off preservation in inverse multi-objective convex optimization. *European Journal of Operational Research*, 270(1), 25–39.
- Chan, T. C., Lee, T., & Terekhov, D. (2019). Inverse optimization: Closed-form solutions, geometry, and goodness of fit. *Management Science*, 65(3), 1115–1135.
- Chang, H. S., Fu, M. C., Hu, J., & Marcus, S. I. (2005). An adaptive sampling algorithm for solving Markov decision processes. *Operations Research*, 53(1), 126–139.

- Chow, J. Y., Ritchie, S. G., & Jeong, K. (2014). Nonlinear inverse optimization for parameter estimation of commodity-vehicle-decoupled freight assignment. *Transportation Research Part E: Logistics and Transportation Review*, 67, 71–91.
- Cozad, A., Sahinidis, N. V., & Miller, D. C. (2014). Learning surrogate models for simulation-based optimization. *AIChE Journal*, 60(6), 2211–2227.
- Domingo, C., Gavaldà, R., & Watanabe, O. (2002). Adaptive sampling methods for scaling up knowledge discovery algorithms. *Data Mining and Knowledge Discovery*, 6(2), 131–152.
- Dong, C., Chen, Y., & Zeng, B. (2018). Generalized inverse optimization through online learning. *Advances in Neural Information Processing Systems, 2018-Decem*(NeurIPS), 86–95.
- Dunning, I., Huchette, J., & Lubin, M. (2017). JuMP: A modeling language for mathematical optimization. *SIAM Review*, 59(2), 295–320.
- Gurobi Optimization, L. (2020). Gurobi optimizer reference manual.
- Hempel, A. B., Goulart, P. J., & Lygeros, J. (2015). Inverse parametric optimization with an application to hybrid system control. *IEEE Transactions on Automatic Control*, 60(4), 1064–1069.
- Heuberger, C. (2004). Inverse combinatorial optimization: A survey on problems, methods, and results. *Journal of Combinatorial Optimization*, 8(3), 329–361.
- Iyengar, G. & Kang, W. (2005). Inverse conic programming with applications. *Operations Research Letters*, 33(3), 319–330.
- Keshavarz, A., Wang, Y., & Boyd, S. (2011). Imputing a convex objective function. *IEEE International Symposium on Intelligent Control - Proceedings*, 613–619.
- Liu, L. & Zhang, J. (2006). Inverse maximum flow problems under the weighted Hamming distance. *Journal of Combinatorial Optimization*, 12(4), 394–407.
- Mohajerin Esfahani, P., Shafieezadeh-Abadeh, S., Hanasusanto, G. A., & Kuhn, D. (2018). Data-driven inverse optimization with imperfect information. *Mathematical Programming*, 167(1), 191–234.
- Saez-Gallego, J. & Morales, J. M. (2018). Short-term forecasting of price-responsive loads using inverse optimization. *IEEE Transactions on Smart Grid*, 9(5), 4805–4814.
- Saez-Gallego, J., Morales, J. M., Zugno, M., & Madsen, H. (2016). A Data-Driven Bidding Model for a Cluster of Price-Responsive Consumers of Electricity. *IEEE Transactions on Power Systems*, 31(6), 5001–5011.
- Sahinidis, N. V. (1996). BARON: A general purpose global optimization software package. *Journal of Global Optimization*, 8(2), 201–205.
- Sahinidis, N. V., Grossmann, I. E., Fornari, R. E., & Chathrathi, M. (1989). Optimization model for long range planning in the chemical industry. *Computers and Chemical Engineering*, 13(9), 1049–1063.

- Schaefer, A. J. (2009). Inverse integer programming. *Optimization Letters*, 3(4), 483–489.
- Schoemaker, P. J. (1991). The quest for optimality: A positive heuristic of science? *Behavioral and Brain Sciences*, 14(2), 205–215.
- Shahmoradi, Z. & Lee, T. (2019). Quantile inverse optimization: Improving stability in inverse linear programming. *arXiv preprint arXiv:1908.02376*.
- Terekhov, A. V., Pesin, Y. B., Niu, X., Latash, M. L., & Zatsiorsky, V. M. (2010). An analytical approach to the problem of inverse optimization with additive objective functions: An application to human prehension. *Journal of Mathematical Biology*, 61(3), 423–453.
- Troutt, M. D., Pang, W. K., & Hou, S. H. (2006). Behavioral estimation of mathematical programming objective function coefficients. *Management Science*, 52(3), 422–434.
- Wang, L. (2009). Cutting plane algorithms for the inverse mixed integer linear programming problem. *Operations Research Letters*, 37(2), 114–116.
- Westermann, K., Lin, J. F.-S., & Kulić, D. (2020). Inverse optimal control with time-varying objectives: application to human jumping movement analysis. *Scientific reports*, 10(1), 1–15.
- Yang, C., Zhang, J., & Ma, Z. (1997). Inverse maximum flow and minimum cut problems. *Optimization*, 40(2), 147–170.
- Zhang, J. & Cai, M. C. (1998). Inverse problem of minimum cuts. *Mathematical Methods of Operations Research*, 47(1), 51–58.
- Zhang, J. & Liu, Z. (1999). A further study on inverse linear programming problems. *Journal of Computational and Applied Mathematics*, 106(2), 345–359.
- Zhang, J. & Xu, C. (2010). Inverse optimization for linearly constrained convex separable programming problems. *European Journal of Operational Research*, 200(3), 671–679.
- Zhang, Q., Grossmann, I. E., & Lima, R. M. (2016). On the relation between flexibility analysis and robust optimization for linear systems. *AIChE Journal*, 62(9), 3109–3123.

Appendix

A-1 Algorithm Pseudocode

Algorithm 4 An efficient algorithm for solving (IOP) in an online environment

```

1: initialize:  $\hat{\mathcal{C}} \leftarrow \mathbb{R}^n$ 
2: for all  $\ell \in \{1, \dots, N\}$  do
3:   solve (P1) $_{\ell}$ , obtain  $\hat{x}_{\ell}^*$ 
4:    $\hat{\mathcal{C}} \leftarrow \hat{\mathcal{C}} \cap \text{cone} \left( \{-a_{\ell t}\}_{t \in \mathcal{T}(\hat{x}_{\ell}^*)} \right)$ 
5:   solve (FP) $_{\ell}$ 
6:   if (FP) $_{\ell}$  is feasible then
7:     solve (P2) with  $\mathcal{I} = \{1, \dots, \ell\}$ , obtain  $\hat{c}_{\ell}^*$ 
8:   else
9:     use Algorithm 1 to solve (P1) $_{[\ell]}$ , warm-start with  $(\hat{x}_1^*, \dots, \hat{x}_{\ell-1}^*)$ 
10:     $\hat{\mathcal{C}} \leftarrow \bigcap_{i \in \mathcal{I}} \text{cone} \left( \{-a_{it}\}_{t \in \mathcal{T}(\hat{x}_i^*)} \right)$ 
11:    solve (P2) with  $\mathcal{I} = \{1, \dots, \ell\}$ , obtain  $\hat{c}_{\ell}^*$ 
12:  end if
13: end for
14:  $\hat{c}^* \leftarrow \hat{c}_{|\mathcal{I}|}^*$ 

```

A-2 Reformulation of (ASP)

In this section, we present reformulations of (ASP) that can be solved using off-the-shelf solvers. We start by reformulating (ASP) as follows:

$$\begin{aligned} & \underset{\eta, y, \gamma, z, A_{\ell}, b_{\ell}, \lambda, \hat{x}_{\ell}, s}{\text{maximize}} && \eta && (14a) \end{aligned}$$

$$\text{subject to} \quad y = \sum_{t \in \mathcal{T}(\hat{x}_i)} \gamma_{it} \left(-\frac{a_{it}}{\|a_{it}\|_2} \right) \quad \forall i \in \mathcal{I} \quad (14b)$$

$$0 \leq \gamma_{it} \leq \epsilon \quad \forall i \in \mathcal{I}, t \in \mathcal{T}(\hat{x}_i) \quad (14c)$$

$$(A_{\ell}, b_{\ell}) \in \Pi \quad (14d)$$

$$\tilde{c}_{\ell-1} + A_{\ell}^{\top} \lambda = 0 \quad (14e)$$

$$A_{\ell} \hat{x}_{\ell} + s = b_{\ell} \quad (14f)$$

$$\lambda \leq Mz \quad (14g)$$

$$s \leq M(e - z) \quad (14h)$$

$$e^{\top} z \geq n \quad (14i)$$

$$\hat{x}_{\ell} \in \mathbb{R}^n, y \in \mathbb{R}^n \quad (14j)$$

$$s_k \geq 0, \lambda_k \geq 0, z_k \geq 0 \quad \forall k \in \mathcal{K} \quad (14k)$$

$$\eta = \min_{\tilde{\gamma}, \delta} e^{\top} \delta \quad (14l)$$

$$\begin{aligned}
\text{subject to} \quad & -\delta \leq y - \sum_{k \in \mathcal{K}} \bar{\gamma}_k z_k a_{\ell k} \\
& y - \sum_{k \in \mathcal{K}} \bar{\gamma}_k z_k a_{\ell k} \leq \delta \\
& \delta \in \mathbb{R}_+^n \\
& \bar{\gamma}_k \geq 0 \quad \forall k \in \mathcal{K},
\end{aligned} \tag{14m}$$

where constraints (14f) - (14k) are the KKT optimality conditions of the linear optimization problem embedded in constraint set (ASPe). Here, we exploit the fact that the binary vector z occurs naturally in the linearized version of the KKT conditions. We also linearize the optimization problem embedded in the objective function of (ASP). Since the sole source of nonlinearity in this inner problem is the norm operator, we introduce a variable δ to linearize it. The rest of the constraints remain the same as in (ASP).

Now, since (14) is a bilevel optimization problem with an LP embedded in it, the problem can be reformulated by replacing the lower-level problem with its optimality conditions, based on strong duality or KKT. Here, we present both reformulations:

1. Reformulating (14) into a single-level problem by using strong duality for the lower-level problem:

$$\text{maximize}_{\eta, y, \gamma, \bar{\gamma}, z, A_\ell, b_\ell, \lambda, \hat{x}, s, \delta, \nu, \omega} \eta \tag{15a}$$

$$\text{subject to} \quad (14b) - (14k) \tag{15b}$$

$$\eta = (\nu - \omega)^\top y \tag{15c}$$

$$-\delta \leq y - \sum_{k \in \mathcal{K}} \bar{\gamma}_k z_k a_{\ell k} \tag{15d}$$

$$y - \sum_{k \in \mathcal{K}} \bar{\gamma}_k z_k a_{\ell k} \leq \delta \tag{15e}$$

$$\nu + \omega \leq e \tag{15f}$$

$$(\nu - \omega)^\top a_{\ell k} z_k \leq 0 \quad \forall k \in \mathcal{K} \tag{15g}$$

$$\delta \in \mathbb{R}_+^n, \bar{\gamma}_k \geq 0 \quad \forall k \in \mathcal{K} \tag{15h}$$

$$\nu \in \mathbb{R}_+^n, \omega \in \mathbb{R}_+^n, \tag{15i}$$

where ν and ω are the dual variables for the lower-level problem in (14). Constraint (15c) imposes the strong duality condition by equating primal and dual objective functions of the lower-level problem in (14). Constraints (15d)-(15e) are the primal feasibility conditions, (15f)-(15g) are the dual feasibility conditions.

2. Reformulating (14) into a single-level problem by using KKT-based optimality conditions for the lower-level problem:

$$\text{maximize}_{\eta, y, \gamma, \bar{\gamma}, z, A_\ell, b_\ell, \lambda, \hat{x}, s, \delta, \nu, \omega, \psi, \phi, s^1, s^2, z^1, z^2, z^3, z^4} \eta \tag{16a}$$

$$\text{subject to} \quad (14b) - (14k) \tag{16b}$$

$$e - \nu - \omega - \psi = 0 \quad (16c)$$

$$(\nu - \omega)^\top a_{\ell k} z_k - \phi_k = 0 \quad \forall k \in \mathcal{K} \quad (16d)$$

$$-\delta + s^1 = y - \sum_{k \in \mathcal{K}} \bar{\gamma}_k z_k a_{\ell k} \quad (16e)$$

$$y - \sum_{k \in \mathcal{K}_M} \bar{\gamma}_k z_k a_{\ell k} + s^2 = \delta \quad (16f)$$

$$s^1 \leq M z^1 \quad (16g)$$

$$\nu \leq M(e - z^1) \quad (16h)$$

$$s^2 \leq M z^2 \quad (16i)$$

$$\omega \leq M(e - z^2) \quad (16j)$$

$$\delta \leq M z^3 \quad (16k)$$

$$\psi \leq M(e - z^3) \quad (16l)$$

$$\bar{\gamma} \leq M z^4 \quad (16m)$$

$$\phi \leq M(e - z^4) \quad (16n)$$

$$\delta \in \mathbb{R}_+^n, \nu \in \mathbb{R}_+^n, \omega \in \mathbb{R}_+^n, \psi \in \mathbb{R}_+^n, s^1 \in \mathbb{R}_+^n, s^2 \in \mathbb{R}_+^n \quad (16o)$$

$$z^1 \in \{0, 1\}^n, z^2 \in \{0, 1\}^n, z^3 \in \{0, 1\}^n \quad (16p)$$

$$\bar{\gamma}_k \geq 0, \phi_k \geq 0, z_k^4 \in \{0, 1\} \quad \forall k \in \mathcal{K}, \quad (16q)$$

where ν , ω , ψ , and ϕ are the Lagrange multipliers of the constraints (14m) of the lower-level problem in (14). Constraints (16c) - (16d) are the stationarity conditions, (16e) - (16f) are the primal feasibility conditions, and (16g) - (16n) represent a linearized version of the complementary conditions corresponding to the constraints (14m) of the lower-level problem in (14).

In our computational case studies, we find that the KKT-based reformulation (16) results in a relatively loose formulation owing to a large number of big-M parameters. Therefore, we employ the strong duality-based (15) to choose input parameters for the case studies in Section 6.

A-3 Additional Results for the Case Study from Section 6.1

In this section, we present some additional results for the customer preference learning case study considered in Section 6.1. Specifically, in Figure 10, we show for each case the evolution of the prediction error for each of the ten random instances of (IOP). Here, one can notice that in case of a failed run with the decomposition algorithm, the solution method is still able to process more than 50% of the experiments. Unlike Algorithm 1, a failed run with Algorithm 2 yields a partial estimate which can be used to generate predictions, albeit with a slightly higher prediction error compared to the mean values observed after the completion of 100 experiments. Further, Figure 10 shows that in the majority of the cases, the prediction error shows a continuous decrease as data from more experiments are added; the spikes observed in the mean curves are a result of just 1 or 2 instances showing a momentary increase in the prediction error. This shows that our IOP formulation is fairly robust in generating accurate predictions even in high-noise situations.

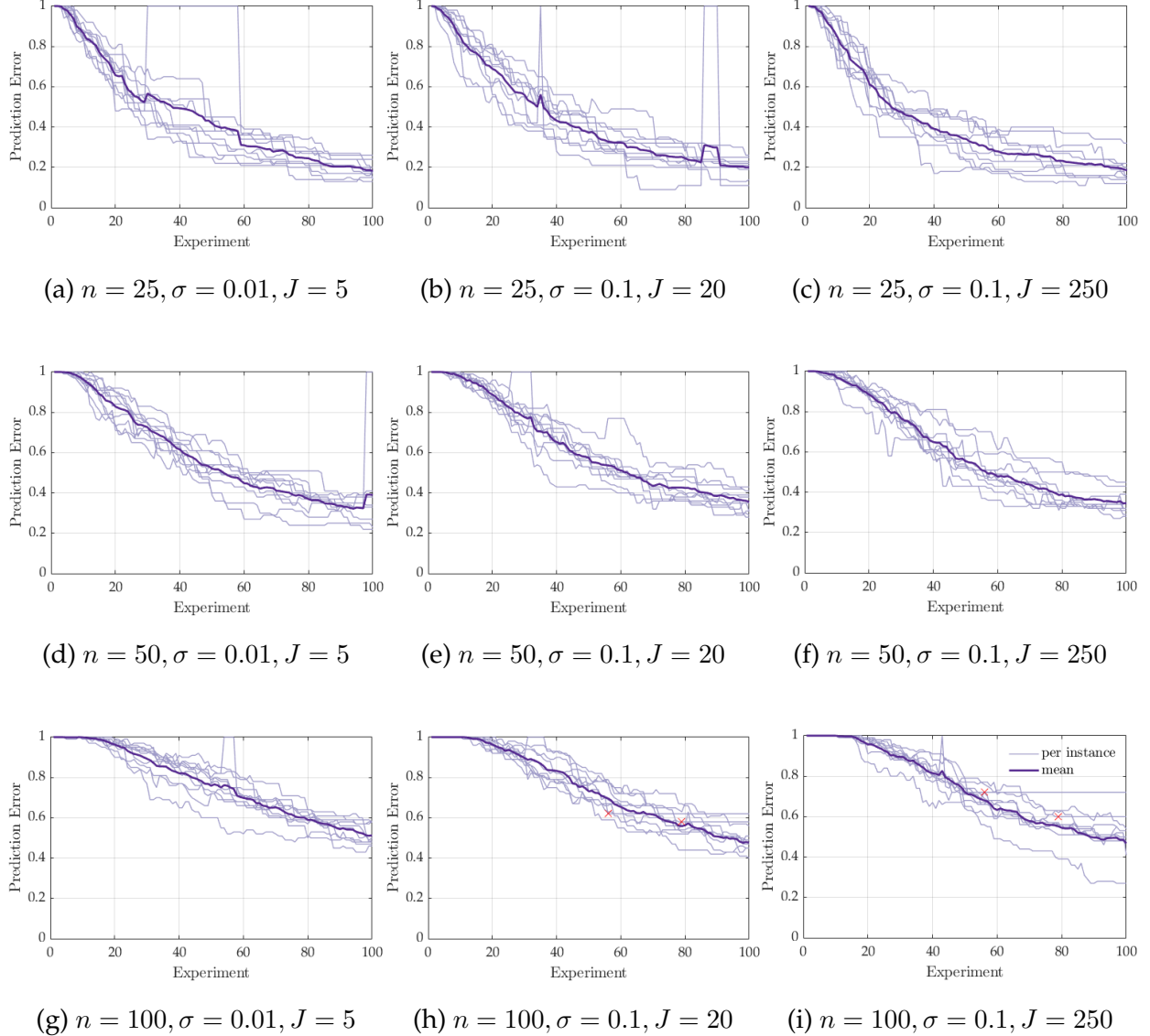


Figure 10: Change in prediction error as experiments are added. Thick lines show mean values across all solved instances, thinner lines show the change in prediction error for all the ten instances. Red cross markers highlight the points at which Algorithm 4 timed out for a specific instance.

A-4 Data for Production Planning Case Study from Section 6.2

We have $q_m^{\min} = 0, q_m^{\max} = 200, x_{pmh}^{\max} = 300$ for all $p \in \mathcal{P}, m \in \mathcal{M}, h \in \mathcal{H}$, and

$$w_{mh}^{\max} = \begin{cases} 300, & \text{if } m \in \{1, \dots, 9\} \\ 0, & \text{otherwise} \end{cases} \quad \forall h \in \mathcal{H}.$$

We assume that the nominal value of demand d_{mh} remains the same across all time periods, i.e. nominal $d_{mh} = D_m \forall h \in \mathcal{H}$. Table 2 lists the values of D_m, q_m^0 used for all products. Materials not listed are raw materials, byproducts, or intermediate products for which D_m and q_m^0 are set to

zero. Table 3 shows the nominal μ_{pm} parameters.

	Materials															
	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
D_m	40	150	30	75	60	30	30	50	150	70	30	75	100	75	250	50
q_m^0	10	5	0	5	5	7	10	2	3	5	10	7	5	5	3	3

Table 2: Values of D_m and q_m^0 parameters used for the case study presented in Section 6.2.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28								
1	-0.58										1																									
2						-0.63						1																								
3	-0.055	-1.25				-0.64					1																									
4		-0.4	-0.69										1																							
5														1																						
6		-0.74												-2.3																						
7													1																							
8			-1												-1.1																					
9																1																				
10													-1.57			-1.26	1																			
11			-1.01														1																			
12			-0.76	-0.28			1																													
13							-1.14																													
14				-0.78									1																							
15																																				
16				-0.6									1																							
17				-0.67																																
18																																				
19																																				
20				-0.35																																
21																																				
22				-0.88	1																															
23					-0.92																															
24				-0.39																																
25					1																															
26				-0.3																																
27																																				
28																																				
29																																				
30				-1.17																																
31					-0.75																															
32				-0.53																																
33																																				
34																																				
35																																				
36																																				
37																																				
38																																				

Table 3: Nominal values of the conversion factors μ_{pm} .