

Using Neural Networks to Solve Linear Bilevel Problems with Unknown Lower Level

IOANA MOLAN, MARTIN SCHMIDT

ABSTRACT. Bilevel problems are used to model the interaction between two decision makers in which the lower-level problem, the so-called follower's problem, appears as a constraint in the upper-level problem of the so-called leader. One issue in many practical situations is that the follower's problem is not explicitly known by the leader. For such bilevel problems with unknown lower-level model we propose the use of neural networks to learn the follower's optimal response for given decisions of the leader based on available historical data of pairs of leader and follower decisions. Integrating the resulting neural network in a single-level reformulation of the bilevel problem leads to a challenging model with a black-box constraint. We exploit Lipschitz optimization techniques from the literature to solve this reformulation and illustrate the applicability of the proposed method with some preliminary case studies using academic and linear bilevel instances.

1. INTRODUCTION

Bilevel optimization has been a highly active field of research in the last decades and has gained increasing attention over the last years. The main reason is that this class of optimization problems can serve as a powerful modeling tool in situations in which one has to model hierarchical decision making; see, e.g., the recent surveys by Beck et al. (2022) and Kleinert et al. (2021) as well as the annotated bibliography by Dempe (2020) to get an overview of the many applications of bilevel optimization. However, this ability also renders bilevel optimization problems very hard to solve both in theory and practice. For instance, even linear bilevel problems are strongly NP-hard (Hansen et al. 1992).

As discussed in the survey by Beck et al. (2022), bilevel optimization problems can also be subject to uncertainty. Moreover, the sources of potential uncertainties are even richer compared to usual, i.e., single-level, optimization. The reason is that not only the problem's data can be uncertain but also the (observation of the) decisions of the two players can be noisy or uncertain. For instance, it might be the case that the leader is not sure about whether the follower can solve the respective lower-level problem to global optimality and might want to hedge against possibly occurring near-optimal solutions; see, e.g., Besançon et al. (2019, 2021). In this short paper, we go even one step further and assume that the leader has no knowledge about an explicit formulation of the lower-level problem. Hence, the leader needs to solve a bilevel optimization problem and the lower level is unknown. What might sound impossible at a first glance can be done, at least approximately, if the bilevel game is repeatedly played so that past pairs of leader decisions and respective responses of the follower can be observed and collected. The obtained set of pairs of decisions can then be used as training data to train a neural network that learns the best-response function of the follower.

Date: October 27, 2022.

2020 Mathematics Subject Classification. 90Cxx, 90C90, 91A26, 91A65.

Key words and phrases. Bilevel optimization, Unknown follower problems, Neural networks, Lipschitz optimization.

The downside of this approach is that the input-output mapping of the neural network can again not be stated using a closed-form expression. This leads to the field of optimization under black-box constraints (Rios and Sahinidis 2013). Fortunately, there is some recent research on deep neural networks with specific activation functions that shows that the input-output mapping of the neural network is Lipschitz continuous and that Lipschitz constants can actually be computed using specifically chosen semidefinite programs (Fazlyab et al. 2019; Pauli et al. 2022). This paves the way for applying the Lipschitz optimization method proposed in Schmidt et al. (2019), see Schmidt et al. (2021) as well, in order to solve the single-level reformulation of the bilevel problem with unknown lower level in which the optimal response of the follower is replaced by the mapping that describes the input-output behavior of the trained neural network. In the Lipschitz-based approach, the only nonlinearity, which is the one given by the input-output mapping of the neural network, is outer-approximated using the Lipschitz constant of the mapping. This outer approximation is tightened from iteration to iteration, leading to a union of polyhedra that form a relaxation of the graph of the nonlinearity, which enables the application of powerful state-of-the-art mixed-integer solvers.

The main contribution of this short note is the combination of the recent results about the Lipschitz continuity of special neural networks with recent publications on Lipschitz optimization. By carrying this out in a careful way and by slightly adapting the method proposed in Schmidt et al. (2019), we obtain a convergent algorithm for this highly challenging situation. We further illustrate the applicability of our approach in a case study based on academic bilevel instances from the literature.

Although there have been some applications of bilevel optimization problems for machine learning (see, e.g., Table 1 in Khanduri et al. (2021)), this is, to the best of our knowledge, the first paper that uses neural networks to solve general bilevel optimization problems. The only other paper we are aware of following this idea is the one by Vlah et al. (2022), who apply deep convolutional neural networks to classic bilevel bidding problems in power markets. However, they focus on the specific application and the specific type of network and its training. In contrast, we focus on the general mathematical idea of replacing unknown lower levels with neural networks. Furthermore, the work by Borrero et al. (2022) presents a sequential learning method for linear bilevel problems under incomplete information. Provided that the strategic players interact with each other multiple times, the authors develop feedback mechanisms to update the missing information for the lower-level objective. Nevertheless, this way to address uncertainty and learning clearly differs from our neural-network based approach. The main idea of this short paper is to give a proof of concept for the proposed method and we, thus, make some simplifications such as that we only consider linear bilevel problems with a scalar upper-level variable and without coupling constraints. We discuss these assumptions in more detail in Section 2 and how the setting can be generalized in Section 6. Let us finally comment on that our approach is related to other methods for solving bilevel optimization problems that rely on using optimal-value or best-response functions (see, e.g., Lozano and Smith (2017) and the references therein). However, our approach is different since we do not work with these functions themselves but with surrogates that we obtain from training a neural network.

The remainder of the paper is structured as follows. In Section 2, we introduce the class of bilevel problems that we consider and pose the main assumptions that are required in what follows. Afterward, in Section 3, we then review the recent literature about computing Lipschitz constants for deep neural networks, which is a prerequisite for the overall solution approach discussed in Section 4. Section 5 contains the case studies that illustrate the applicability of our approach for small

instances from the literature. Finally, we conclude in Section 6 and discuss some topics for future research.

2. PROBLEM STATEMENT

We consider linear bilevel problems of the general form

$$\min_{x \in X, y} c^\top x + d^\top y \quad (1a)$$

$$\text{s.t. } Ax \geq a, \quad (1b)$$

$$\underline{x} \leq x \leq \bar{x}, \quad x \in \mathbb{R}^{n_x}, \quad (1c)$$

$$y \in \Psi(x), \quad (1d)$$

where $\Psi(x)$ is the set of optimal solutions of the x -parameterized problem

$$\min_{y \in Y} f^\top y \quad (2a)$$

$$\text{s.t. } Cx + Dy \geq b, \quad (2b)$$

with $c \in \mathbb{R}^{n_x}$, $d, f \in \mathbb{R}^{n_y}$, $A \in \mathbb{R}^{m \times n_x}$, $C \in \mathbb{R}^{\ell \times n_x}$, $D \in \mathbb{R}^{\ell \times n_y}$, $a \in \mathbb{R}^m$, and $b \in \mathbb{R}^\ell$. Problem (1) is called the upper-level or leader's problem. The decision variables of the leader are $x \in \mathbb{R}^{n_x}$. Problem (2) is called the lower-level or follower's problem and has the decision variables $y \in \mathbb{R}^{n_y}$. Note that we consider the optimistic bilevel problem. Hence, the leader also optimizes over y if the lower-level problem is not uniquely solvable. The set $\Omega := \{(x, y) \in X \times Y : (1b), (1c), (2b)\}$ is called the shared constraint set and the set $\mathcal{F} := \{(x, y) \in \Omega : y \in \Psi(x)\}$ is called the bilevel feasible set. With \mathcal{F}_x , we denote its projection onto the x variables. For what follows, we make the ensuing assumptions.

- Assumption 1.** (i) *The upper-level decision $x \in \mathbb{R}^{n_x}$ is scalar, i.e., $n_x = 1$.*
(ii) *For all upper-level decisions x for which $\Psi(x)$ is non-empty, we have $|\Psi(x)| = 1$. Hence, if the lower level is feasible and bounded, then it has a unique optimal solution y . Consequently, we can write $y = \Psi(x)$ instead of $y \in \Psi(x)$.*
(iii) *The solution-set mapping $\Psi(\cdot)$ is Lipschitz continuous.*

According to Assumption 1, $\Psi(x)$ is a singleton, leading to a one-to-one correspondence between follower's optimal responses and the upper-level decisions. Moreover, the mapping $\Psi(\cdot)$ is polyhedral, i.e., its graph is a finite union of polyhedra; see Theorem 3.1 in Dempe (2002). Furthermore, the bilevel feasible set is connected because we have no coupling constraints and, thus, $\Psi_i(x)$ is a Lipschitz continuous and piecewise linear function in x for all $i \in [n_y] := \{1, \dots, n_y\}$.

Lower-level uniqueness is particularly important when it comes to training a neural network to learn the optimal response of the follower for a given upper-level decision. It ensures that, during supervised learning, there is only a single output y to be learned for a given input x . The assumption of a scalar leader's decision can be generalized and is mainly taken for the ease of presentation. We will discuss this in more detail in our conclusion in Section 6.

3. USING NEURAL NETWORKS TO APPROXIMATE THE FOLLOWER'S RESPONSE

Stating the bilevel problem (1) relies on the knowledge about an explicit formulation for the lower-level problem (2). In the case in which the follower's problem (2) is not known by the leader but past pairs (x, y) of leader and follower decisions are available, we propose using this historical data to learn the optimal response $y = \Psi(x)$ of the follower using neural networks.

Such (x, y) -pairs naturally arise in many applications. For instance, for pricing models, the upper-level variable x is the price set by the leader for a certain good

and y is the amount of this good bought by the follower. Both the price and the bought goods can be observed and collected to obtain (x, y) -pairs to train a neural network. Similar situations appear in many other fields such as in bilevel optimization for market design, transportation, or security. Obviously, it is required that the game modeled by the bilevel problem is played repeatedly so that large enough training sets can be collected over time.

In what follows, we use a neural network to learn $\Psi(\cdot)$, which will then replace the lower level and turn the bilevel model into the single-level problem

$$\min_{x, y} \quad cx + d^\top y \quad (3a)$$

$$\text{s.t.} \quad Ax \geq a, \quad \underline{x} \leq x \leq \bar{x}, \quad x \in \mathbb{R}, \quad (3b)$$

$$g_i(x) = y_i, \quad i \in [n_y], \quad (3c)$$

where g_i is the function corresponding to the neural network for the i th follower's response y_i , $i \in [n_y]$. Thus, we assume $g_i(x) \approx \Psi_i(x)$. In what follows, we exploit some recent results from the literature to show that g_i is Lipschitz continuous and that Lipschitz constants can indeed be computed. This property of the used neural networks is vital for the decomposition method we use to solve Problem (3); see Section 4.

3.1. Lipschitz Constants of Neural Networks. In this paper, we use the LipSDP-Neuron method published in Fazlyab et al. (2019) to compute Lipschitz constants of the neural network functions g_i , $i \in [n_y]$. That is, we compute a constant $L_i \geq 0$ that satisfies

$$|g_i(x_1) - g_i(x_2)| \leq L_i |x_1 - x_2| \quad \text{for all } x_1, x_2 \in \mathbb{R}$$

and for all $i \in [n_y]$. The main idea in Fazlyab et al. (2019) is to replace the nonlinear activation functions at the nodes of a neural network by so-called incremental quadratic constraints, which then allows to state the problem of estimating Lipschitz constants as a semidefinite program (SDP). It is worth noting that the most complex and, hence, the most accurate version of LipSDP in Fazlyab et al. (2019) is shown to be wrong in Pauli et al. (2022). Thus, we use the second of the three versions of LipSDP, namely LipSDP-Neuron.

We now describe the quadratic constraints that we use to replace all activation functions $\phi(x) = [\varphi(x_1), \dots, \varphi(x_n)]^\top : \mathbb{R}^n \rightarrow \mathbb{R}^n$ in a network layer, where the same slope-restricted function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is applied to each component of ϕ . Here and in what follows, we call a function φ slope-restricted w.r.t. $0 \leq \alpha < \beta < \infty$ if

$$\alpha \leq \frac{\varphi(y) - \varphi(x)}{y - x} \leq \beta \quad (4)$$

holds for all $x, y \in \mathbb{R}$. This definition states that the slope of the line connecting any two points x and y on the graph of φ is bounded by α and β . It is easy to see that the Rectified Linear Unit (ReLU) activation function defined as $\varphi(x) := \max\{0, x\}$ is slope-restricted with $\alpha = 0$ and $\beta = 1$; see Goodfellow et al. (2016). Furthermore, if the activation function φ is slope-restricted w.r.t. $[\alpha, \beta] = [0, 1]$, then so is the vector-valued function $\phi(x) = [\varphi(x_1), \dots, \varphi(x_n)]$, which contains all activation functions in a network layer (Fazlyab et al. 2022), if we apply the definition in (4) component-wise.

The following lemma shows that the slope property (4) can be written as a quadratic constraint. We use the notation \mathbb{S}^n for the set of all symmetric $n \times n$ matrices.

Lemma 1 (Based on Lemma 1 in Fazlyab et al. (2019)). *Suppose $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is slope-restricted w.r.t. α and β . Moreover, we define the set*

$$\mathcal{T}_n := \left\{ T \in \mathbb{S}^n : T = \sum_{i=1}^n \lambda_{ii} e_i e_i^\top, \lambda_{ii} \geq 0 \right\} \quad (5)$$

of diagonal matrices T with nonnegative entries. Then, for any $T \in \mathcal{T}_n$, the function $\phi(x) = [\varphi(x_1), \dots, \varphi(x_n)]^\top : \mathbb{R}^n \rightarrow \mathbb{R}^n$ satisfies the quadratic constraint

$$\begin{pmatrix} x - y \\ \phi(x) - \phi(y) \end{pmatrix}^\top \begin{bmatrix} -2\alpha\beta T & (\alpha + \beta)T \\ (\alpha + \beta)T & -2T \end{bmatrix} \begin{pmatrix} x - y \\ \phi(x) - \phi(y) \end{pmatrix} \geq 0 \quad (6)$$

for all $x, y \in \mathbb{R}^n$.

The statement of the lemma above is based on Lemma 1 in Fazlyab et al. (2019) and has been modified according to the corrections published in Pauli et al. (2022). There, the authors give a counterexample that shows that the original lemma in Fazlyab et al. (2019) is wrong. To be more specific, they show that there exists a matrix T built according to the definition of the set \mathcal{T}_n given in Fazlyab et al. (2019) that violates (6).

Assuming that all activation functions $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ are the same, a feed-forward neural network $f(x) : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{\ell+1}$ can be written compactly as

$$B\mathbf{x} = \phi(A\mathbf{x} + b) \quad \text{and} \quad f(x) = C\mathbf{x} + b^\ell, \quad (7)$$

where $\mathbf{x} = [(x^0)^\top, \dots, (x^\ell)^\top]^\top$ is the concatenation of the input values at every layer of the network, ℓ is the number of layers, $x^i \in \mathbb{R}^{n_i}$ for all $i \in \{1, \dots, \ell\}$, ϕ is the vector-valued function which applies the activation function φ to every entry of the input vector, and

$$A = \begin{bmatrix} W^0 & 0 & \dots & 0 & 0 \\ 0 & W^1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & W^{\ell-1} & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & I_{n_1} & 0 & \dots & 0 \\ 0 & 0 & I_{n_2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & I_{n_\ell} \end{bmatrix}, \quad (8)$$

$$C = [0 \quad \dots \quad 0 \quad W^\ell], \quad b = [(b^0)^\top, \dots, (b^{\ell-1})^\top]^\top$$

holds, where W^i is the weight matrix connecting layer i with layer $i + 1$, and I_{n_i} is the $n_i \times n_i$ identity matrix. The next theorem is the central result in Fazlyab et al. (2019) and shows that the Lipschitz constant of a neural network is the solution of an SDP in which the matrix T defined in Lemma 1 serves as a decision variable.

Theorem 2 (Theorem 2 in Fazlyab et al. (2019)). *Consider an ℓ -layer and fully connected neural network given by (7). Let $n = \sum_{k=1}^\ell n_k$ be the total number of hidden neurons and suppose that the activation functions are slope-restricted w.r.t. α and β . Define \mathcal{T}_n as in (5), A and B as in (8), and consider the matrix inequality*

$$M(\rho, T) = \begin{bmatrix} A \\ B \end{bmatrix}^\top \begin{bmatrix} -2\alpha\beta T & (\alpha + \beta)T \\ (\alpha + \beta)T & -2T \end{bmatrix} \begin{bmatrix} A \\ B \end{bmatrix} + R \preceq 0 \quad (9)$$

with

$$R = \begin{bmatrix} -\rho I_{n_0} & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & (W^\ell)^\top W^\ell \end{bmatrix}.$$

If (9) holds for some $(\rho, T) \in \mathbb{R}_{\geq 0} \times \mathcal{T}_n$, then $\|f(x) - f(y)\|_2 \leq \sqrt{\rho} \|x - y\|_2$ holds for all $x, y \in \mathbb{R}^{n_0}$.

As a result of Theorem 2, a Lipschitz constant for multi-layer networks can be computed by solving

$$\min_{\rho, T} \quad \rho \quad \text{s.t.} \quad M(\rho, T) \preceq 0, \quad T \in \mathcal{T}_n, \quad (10)$$

where $(\rho, T) \in \mathbb{R}_+ \times \mathcal{T}_n$ are the decision variables. Furthermore, $M(\rho, T)$ is linear in ρ and T and the set \mathcal{T}_n is convex. Hence, Problem (10) is a semidefinite program.

4. SOLUTION APPROACH

The neural network constraint (3c) turns model (3) into a challenging problem. In particular, for complex and large neural networks, it cannot be expected to get a closed-form expression for g_i , $i \in [n_y]$, that has reasonable properties required for optimization. In any case, the resulting constraints will be nonlinear, nonconvex, and nonsmooth. However, we can evaluate these constraints and we can compute their Lipschitz constants. Thus, it is reasonable to make the following assumption.

Assumption 2. *An oracle is available that evaluates $g_i(\cdot)$ for all $i \in [n_y]$ and all g_i are globally Lipschitz continuous on $\underline{x} \leq x \leq \bar{x}$ with known global Lipschitz constant L_i .*

To solve Problem (3), we use the decomposition method published in Schmidt et al. (2019) but modify it slightly so that we can apply it to our setting. We assume Lipschitz continuity of the problematic nonlinearities and use the corresponding Lipschitz constants to build a mixed-integer linear problem (MILP) that is a relaxation of the original model (3). We refine this relaxation in every iteration until a satisfactory solution is found or until the problem is shown to be infeasible. A satisfactory solution is formally defined to be an ε -feasible point, i.e., a point that solves

$$\min_{x, y} \quad cx + d^\top y \quad (11a)$$

$$\text{s.t.} \quad Ax \geq a, \quad \underline{x} \leq x \leq \bar{x}, \quad x \in \mathbb{R}, \quad (11b)$$

$$|g_i(x) - y_i| \leq \varepsilon, \quad i \in [n_y], \quad (11c)$$

where $\varepsilon \geq 0$ a user-specified tolerance. Note that this relaxation of Problem (3) only affects the nonlinear functions g_i , $i \in [n_y]$.

4.1. Core Ideas and Notation. The main idea of the decomposition method is to relax the graphs of g_i with help of a set Ω_i . The set Ω_i is given by linear constraints built around g_i using the corresponding Lipschitz constant L_i . The resulting relaxation Ω_i is then refined in each iteration k such that Ω_i^k can be written as a union of polytopes

$$\Omega_i^k = \bigcup_{j \in J_i^k} \Omega_i^k(j),$$

which converges towards the graph of g_i for $k \rightarrow \infty$. The union of polytopes is uniquely defined in each iteration k by a set of points on the x -axis,

$$\mathcal{X}_i^k := \{x_i^{k,0}, x_i^{k,1}, \dots, x_i^{k,|J_i^k|}\},$$

where $x_i^{k,j} \in \mathbb{R}$ are scalar values for $j \in \{0\} \cup J_i^k = \{0\} \cup \{1, \dots, |J_i^k|\}$ and

$$\underline{x}_i =: x_i^{k,0} < x_i^{k,1} < \dots < x_i^{k,|J_i^k|} := \bar{x}_i.$$

Note that while x_i are scalar, the index i indicates that the sets \mathcal{X}_i can develop differently from each other, specifically to every $i \in [n_y]$. In other words, the relaxations of g_i can be refined individually in each iteration k . This relaxation Ω_i^k

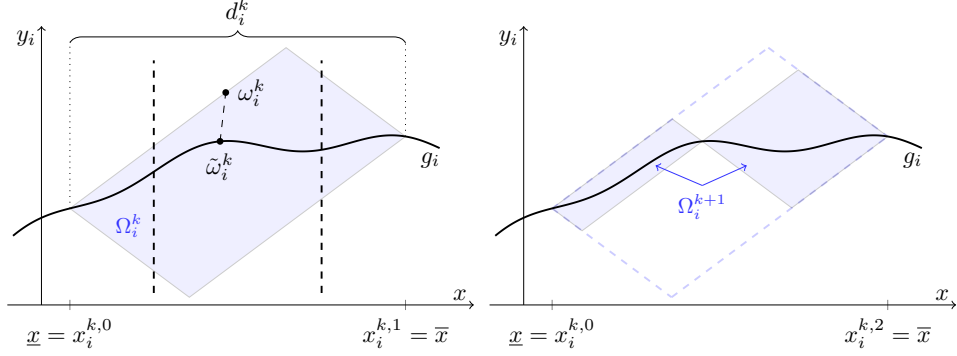


FIGURE 1. Left: The subproblem is solved on the set $\tilde{\Omega}_i^k$, which is bounded by the dashed, vertical lines. Right: The feasible set of the master problem in iteration $k + 1$, after \tilde{x}_i has been added to \mathcal{X}_i^k . The dashed lines depict the old feasible set.

of g_i is improved by adding a new point $x_i^{k,j}$ to \mathcal{X}_i^k . The polytopes of the union mentioned above are given by

$$\begin{aligned} \Omega_i^k(j) = \Big\{ (x, y_i) \in \mathbb{R}^2 : & x_i^{k,j-1} \leq x \leq x_i^{k,j}, \\ & y_i \leq g_i(x_i^{k,j-1}) + L_i(x - x_i^{k,j-1}), \\ & y_i \geq g_i(x_i^{k,j-1}) - L_i(x - x_i^{k,j-1}), \\ & y_i \leq g_i(x_i^{k,j}) + L_i(x_i^{k,j} - x), \\ & y_i \geq g_i(x_i^{k,j}) - L_i(x_i^{k,j} - x) \Big\} \end{aligned} \quad (12)$$

for $j \in J_i^k$. Visually speaking, (12) states that the polytopes $\Omega_i^k(j)$ are all quadrilaterals with two vertices $x_i^{k,j-1}$ and $x_i^{k,j}$ on the graph of g_i ; see Figure 1. To understand how a point is added to \mathcal{X}_i , we discuss the two problems solved in each iteration of the algorithm.

4.2. The Master Problem. The master problem $(M(k))$ is solved in each iteration over the sets Ω_i^k . The problem is formally given by

$$\begin{aligned} \min_{\omega} \quad & cx + d^\top y \\ \text{s.t.} \quad & Ax \geq a, \quad \underline{x} \leq x \leq \bar{x}, \quad x \in \mathbb{R}, \\ & \omega_i \in \Omega_i^k, \quad i \in [n_y], \end{aligned} \quad (M(k))$$

with $\omega = (x, y)$ and $\omega_i := (x, y_i)$. Note that ω_i is not the i th entry in ω , since $\omega \in \mathbb{R}^{1+n_y}$ is the vector containing x and y while $\omega_i \in \mathbb{R}^2$ contains x and y_i . If the solution ω^k with $x^k \in \mathbb{R}$ and $y^k \in \mathbb{R}^{n_y}$ of $(M(k))$ is ε -feasible, then Problem (3) is approximately solved. According to Schmidt et al. (2019), the master

problem $(\mathbf{M}(k))$ can be modeled as the mixed-integer linear problem

$$\min_{\omega} \quad cx + d^\top y \quad (13a)$$

$$\text{s.t.} \quad Ax \geq a, \quad \underline{x} \leq x \leq \bar{x}, \quad x \in \mathbb{R}, \quad (13b)$$

$$-M(1 - z_i^{k,j}) + x_i^{k,j-1} \leq x \leq x_i^{k,j} + M(1 - z_i^{k,j}), \quad i \in [n_y], j \in J_i^k, \quad (13c)$$

$$y_i \leq g_i(x_i^{k,j-1}) + L_i(x - x_i^{k,j-1}) + M(1 - z_i^{k,j}), \quad i \in [n_y], j \in J_i^k, \quad (13d)$$

$$y_i \geq g_i(x_i^{k,j-1}) - L_i(x - x_i^{k,j-1}) - M(1 - z_i^{k,j}), \quad i \in [n_y], j \in J_i^k, \quad (13e)$$

$$y_i \leq g_i(x_i^{k,j}) + L_i(x_i^{k,j} - x) + M(1 - z_i^{k,j}), \quad i \in [n_y], j \in J_i^k, \quad (13f)$$

$$y_i \geq g_i(x_i^{k,j}) - L_i(x_i^{k,j} - x) - M(1 - z_i^{k,j}), \quad i \in [n_y], j \in J_i^k, \quad (13g)$$

$$\sum_{j \in J_i^k} z_i^{k,j} = 1, \quad i \in [n_y], \quad (13h)$$

$$z_i^{k,j} \in \{0, 1\}, \quad i \in [n_y], j \in J_i^k. \quad (13i)$$

The constraint $\omega_i \in \Omega_i^k$, $i \in [n_y]$, in $(\mathbf{M}(k))$ is modeled here using big- M constraints. The binary variables z in (13i) and (13h) ensure that only one polytope is active for all $i \in [n_y]$.

4.3. The Subproblem. On the other hand, if the solution ω^k is not ε -feasible, the relaxation Ω_i^k of the graph of g_i needs to be refined. This is achieved by solving a subproblem to find a point on the graph of g_i , which is as close as possible to the solution $\omega_i^k \in \Omega_i^k$ of the master problem $(\mathbf{M}(k))$. We then use this point to refine the relaxation.

Our subproblem differs from the original method described in Schmidt et al. (2019). There, an optimization problem over nonlinearities and other linear constraints is solved. This is not possible in the setting we consider here, or is at least extremely challenging, due to the nonconvex and nonsmooth nature of neural networks.

The subproblem is solved only for the particular polytope $\Omega_i^k(j_i^k) \subset \Omega_i^k$ with $\omega_i^k \in \Omega_i^k(j_i^k)$. That is, we look for a point on the graph of g_i , which is contained in $\Omega_i^k(j_i^k)$. Additionally, the feasible set of the subproblem is further reduced to only allow for a solution in an inner-approximation of the respective polytope. This way we ensure that newly found points do not accumulate at an already existing value in \mathcal{X}_i^k . For a given $j \in J_i^k$, this subset is defined as

$$\tilde{\Omega}_i^k(j) = \Omega_i^k(j) \cap \hat{\Omega}_i^k(j),$$

with

$$\hat{\Omega}_i^k(j) = \left\{ (x, y_i) \in \mathbb{R}^2 : x_i^{k,j-1} + \frac{1}{4}d_i^{k,j} \leq x \leq x_i^{k,j} - \frac{1}{4}d_i^{k,j} \right\}$$

and $d_i^{k,j} = x_i^{k,j} - x_i^{k,j-1}$ is the length of the corresponding segment on the x -axis. The constant 0.25 can be replaced by any value in $(0, 0.5)$. The left plot in Figure 1 indicates the set $\hat{\Omega}_i^k$ with vertical dashed lines.

To solve the subproblem, we sample points \tilde{x}_i on the x -axis segment corresponding to $\hat{\Omega}_i^k(j_i^k)$ and evaluate g_i at these points. Then, the solution of the subproblem is given by the computed point $\tilde{\omega}_i := (\tilde{x}_i, g_i(\tilde{x}_i))$ that is closest to the solution $\omega_i^k = (x^k, y_i^k)$ of the master problem w.r.t. the Euclidean distance; see Figure 1. In other words, we choose from a finite set of points in $\hat{\Omega}_i^k(j_i^k)$ the one that is on the graph of g_i and as close as possible to the solution ω_i^k of $(\mathbf{M}(k))$. Note that we solve the subproblem only for those $i \in [n_y]$ that are not ε -feasible, i.e., if $|g_i(x^k) - y_i^k| > \varepsilon$ holds.

The solution \tilde{x}_i^k of the subproblem is added to the set \mathcal{X}_i^k and thus refines the relaxation of g_i ; see Figure 1 (right). While x^k in ω_i^k is the solution of $(\mathbf{M}(k))$ and thus shared by all $i \in [n_y]$, the solution \tilde{x}_i^k of the subproblem unique to $i \in [n_y]$.

4.4. Algorithm and Convergence Properties. The Lipschitz decomposition method is formally given in Algorithm 1. There, the master problem (M(k)) is solved in each iteration k and the algorithm checks if the solution ω^k is ε -feasible. If this is not the case for all $i \in [n_y]$, the polytope $\Omega_i^k(j_i^k)$ containing the solution ω_i^k is identified using the index of the variables $z_i^{k,j}$ in the MILP (13) for all ε -infeasible $i \in [n_y]$. Then, a new point $\tilde{\omega}_i^k = (\tilde{x}_i^k, \tilde{y}_i^k)$ is found on a subset $\tilde{\Omega}_i^k(j_i^k) \subset \Omega_i^k(j_i^k)$ and \tilde{x}_i^k is added to \mathcal{X}_i^k , refining the approximation of g_i . Notice that while x is scalar, the index i in \tilde{x}_i^k signals that the new point found on the x -axis in iteration k is specific to the relaxation of function g_i . Moreover, the number of points in \mathcal{X}_i^k is at most k in every iteration, which means that $|J_i^k| \leq k$ for all $i \in [n_y]$.

Algorithm 1 Lipschitz Decomposition Method

Input: Problem (3) and $\varepsilon > 0$

Output: An approximate globally optimal and ε -feasible point for Problem (3) or indication of infeasibility

```

1: Set  $k \leftarrow 1$  and initialize  $\mathcal{X}_i^k = \{\underline{x}, \bar{x}\}$  for all  $i \in [n_y]$ .
2: while true do
3:   Solve the master problem (13) to global optimality.
4:   if (13) is infeasible then
5:     return "Problem (3) is infeasible."
6:   end if
7:   Let  $\omega^k = (x^k, y^k)$  denote the optimal solution of (13).
8:   if  $|g_i(x^k) - y_i^k| \leq \varepsilon$  for all  $i \in [n_y]$  then
9:     return  $\omega^k$ .
10:  end if
11:  Determine the polytopes  $j_i^k \in J_i^k$  for all  $i \in [n_y]$ .
12:  Solve the subproblems for all  $i \in [n_y]$  with  $|g_i(x^k) - y_i^k| > \varepsilon$  and let  $\tilde{\omega}_i^k = (\tilde{x}_i^k, \tilde{y}_i^k)$  denote the optimal solution.
13:  for  $i \in [n_y]$  do
14:    if  $|g_i(x^k) - y_i^k| > \varepsilon$  then
15:      Set  $\mathcal{X}_i^{k+1} \leftarrow \mathcal{X}_i^k \cup \{\tilde{x}_i^k\}$ .
16:    else
17:      Set  $\mathcal{X}_i^{k+1} \leftarrow \{\mathcal{X}_i^k\}$ .
18:    end if
19:  end for
20:  Increase  $k \leftarrow k + 1$ .
21: end while

```

We remark that the new point \tilde{x}_i^k found by the subproblem splits the original quadrilateral $\Omega_i^k(j_i^k)$ into two smaller ones; see Figure 1. For the two new polytopes, we use the notations $\Omega_i^k(j_1^k)$ and $\Omega_i^k(j_2^k)$. Hence, the union of polytopes Ω_i^{k+1} that approximates g_i in iteration $k + 1$ is given by

$$\Omega_i^{k+1} = \Omega_i^k(j_1^k) \cup \Omega_i^k(j_2^k) \cup \bigcup_{j \neq j_1^k, j_2^k \in J_i^k} \Omega_i^k(j).$$

and we have the following lemma.

Lemma 3 (See Lemma 4 in Schmidt et al. (2019)). *There exists a constant $\delta > 0$ depending only on ε and L such that as long as Algorithm 1 does not terminate in Line 5 or 9, there exists a constant $\delta^k > \delta$ for every k with*

$$\text{Vol}(\Omega_i^k(j_1^k)) + \text{Vol}(\Omega_i^k(j_2^k)) = \text{Vol}(\Omega_i^k(j_i^k)) - \delta^k.$$

Theorem 4 (See Theorem 1 in Schmidt et al. (2019)). *There exists a $K < \infty$ such that Algorithm 1 either terminates with an approximate globally optimal point ω^k or with the indication of infeasibility in an iteration $k \leq K$.*

Similarly to Theorem 1 in Schmidt et al. (2019), Theorem 4 follows from the fact that, according to Lemma 3, the volume of Ω_i^k decreases by a positive value δ^k in each iteration that is uniformly bounded away from zero.

However, due to the fact that Ω_i is refined in each iteration to better approximate g_i , the master problem $(M(k))$ grows linearly over the course of the iterations, leading to one additional binary variable $z_i^{k,j}$ and some additional linear constraints. Consequently, the computational effort increases in every iteration.

5. CASE STUDY

In this section, we illustrate the applicability of Algorithm 1 on the basis of a set of exemplary case studies. We first discuss some implementation details and then present the results of the algorithm.

5.1. Implementation Details. We now explain how we generate the (x, y) -pairs for the considered instances. We also discuss the used neural networks and their training, the sampling of points for the subproblem, and how we verify the solutions obtained with our algorithm.

We generate the (x, y) -pairs as follows. First, we solve the high-point relaxation of the bilevel problem but with a different objective function in which we either minimize or maximize the upper-level variable x to get the set \mathcal{F}_x . Then, we equidistantly sample in this interval and solve the parametric lower-level problem for the given values of x , obtain the lower-level problem's solution y , and store the point (x, y) with $y = \Psi(x)$ in our data set. The resulting data set entirely consists of bilevel feasible points. The obtained (x, y) -pairs are then used to train neural networks to learn the optimal follower's response. Note that using the modified high-point relaxation to obtain the x -interval used for sampling is not possible in reality if the lower-level problem is not known. However, a generation procedure would not be required in reality at all since the set of (x, y) -pairs for training the neural networks consist of observed data from the past.

The lower-level problem of every instance is then assumed to be unknown. Instead, for every considered problem, the optimal follower's response for a given leader's decision x , i.e., $\Psi_i(x)$, is learned with feed-forward neural networks with ReLU activation functions at every node, for all $i \in [n_y]$; see, e.g., Goodfellow et al. (2016). The functions learned with ReLU networks are inherently piecewise linear and, thus, Lipschitz continuous. Moreover, for all instances the weights of the network are updated via the Adam optimizer (Ruder 2016). We vary the learning rate, the number of epochs, and the network architecture depending on the instance at hand.

The bilevel feasible set \mathcal{F} and its projection \mathcal{F}_x onto the x -variables also depend on the unknown lower level. Hence, we attempt to deduce \mathcal{F} and \mathcal{F}_x from the available (x, y) -pairs. Due to the fact that the interval $[\underline{x}, \bar{x}]$ with $\underline{x} \leq x \leq \bar{x}$ can be larger than \mathcal{F}_x , we do not use \underline{x} and \bar{x} to initialize Ω_i ; see (12). Instead, we use the closed interval generated by the smallest and the largest x -values in the available training set as a proxy for \mathcal{F}_x .

Given the trained networks g_i , we solve the master problem and, subsequently, the subproblems as described in Section 4. To solve the subproblem, we evaluate the functions g_i at $p = 100$ points on the corresponding x -axis segment. If s points have already been evaluated on this segment in earlier iterations, then only $p - s$ new equidistantly distributed points $(x, g_i(x))$ are computed.

TABLE 1. Basic information about the instance-specific neural networks. All networks are trained for 200 epochs and have training losses close to zero. The learning rates are rounded to three decimals.

Instance	Training set size	Network architecture	Learning rate
Bard (1984)	30	[1, 15, 10, 15, 1]	0.055
Bialas and Karwan (1984)	30	[1, 20, 10, 20, 1]	0.049
Clark and Westerberg (1988)	30	[1, 10, 15, 10, 1]	0.023
Haurie et al. (1990)	30	[1, 10, 15, 10, 1]	0.049
Liu and Hart (1994)	30	[1, 20, 10, 20, 1]	0.094
Moore and Bard (1990)	30	[1, 5, 5, 1]	0.074

TABLE 2. Lipschitz constants as computed with the LipSDP-Neuron method using MOSEK 9.3.20; see Section 3.

Instance	Lipschitz constant	Computation time
Bard (1984)	5.77	0.049
Bialas and Karwan (1984)	7.69	0.052
Clark and Westerberg (1988)	0.98	0.045
Haurie et al. (1990)	7.3	0.059
Liu and Hart (1994)	21.65	0.081
Moore and Bard (1990)	3.02	0.022

Finally, Algorithm 1 stops with the indication that Problem (3) is infeasible or with an ε -feasible solution, where $\varepsilon = 10^{-5}$ is used in our experiments. The solution computed by our algorithm is compared with the solution obtained by solving the mixed-integer KKT reformulation with sufficiently large big- M constants; see, e.g., Kleinert et al. (2021).

We implemented the LipSDP-Neuron method using the `cvxpy` 1.1.13 package and the SDP solver MOSEK 9.3.20. All occurring linear or mixed-integer linear problems (in particular, $(M(k))$) are solved using Gurobi 9.5.1. All neural networks have been trained using the Python library `torch` 1.11.0. All computations have been executed on a Intel® Core™i7-10510U CPU with 8 cores of 1.8 GHz each and 32 GB RAM.

5.2. Discussion of the Results. We apply the Lipschitz decomposition method to 6 instances of linear bilevel problems from the literature. All instances have a scalar upper-level decision variable x so that Assumption 1 is satisfied. For all 6 instances we use 60% of the (x, y) -pairs for training and the rest for validation. Alternative proportions of training and validation set sizes could be used as well if appropriate. Furthermore, given that $\Psi_i(\cdot), i \in [n_y]$, is piecewise linear in our context, rather small data sets are often already enough to find good solutions for the considered instances.

Table 1 shows the training set sizes, the configurations of the neural networks, and the used learning rates. Table 2 contains the corresponding Lipschitz constants computed for the networks described in the previous table as well as the time required to compute them. Finally, Table 3 displays the solutions obtained with Algorithm 1 (as well as the required number of iterations) next to the ones computed by solving the KKT reformulation. All computation times are given in seconds.

TABLE 3. Algorithm 1 solves the instances above using the networks in Table 1 and the corresponding Lipschitz constants in Table 2.

Instance	Solution of KKT reform.	Solution of Algorithm 1	Nr. of iterations	Comp. time
Bard (1984)	(7.2, 1.6)	(7.07, 1.54)	12	0.3
Bialas and Karwan (1984)	(16, 11)	(15.67, 10.34)	55	2.23
Clark and Westerberg (1988)	(19, 14)	(18.63, 13.76)	33	1.42
Haurie et al. (1990)	(12, 3)	(11.79, 3.15)	49	4.44
Liu and Hart (1994)	(4, 4)	(3.92, 3.67)	58	3.22
Moore and Bard (1990)	(0, 1.5)	(0, 1.4999)	31	1.21

5.2.1. *The Performance of Algorithm 1.* As discussed in Section 4.4, the computational effort of solving the master problem naturally increases over the course of the iterations. This can also be seen clearly in Figure 2 for 4 exemplarily chosen instances. Furthermore, the computed Lipschitz constants determine the volume of the polytopes Ω_i^k , $i \in [n_y]$. Thus, according to Lemma 3 and Theorem 4, the convergence of Algorithm 1 directly depends on the constants L_i .

One way to keep Algorithm 1 most efficient is to compute Lipschitz constants L_i that are as small as possible. To illustrate this, we consider the instance (Moore and Bard 1990)

$$\max_{x,y} \quad F(x,y) = -x - 2x \quad (14a)$$

$$\text{s.t.} \quad y \in \Psi(x), \quad (14b)$$

where $\Psi(x)$ is the set of optimal solutions of the x -parameterized linear lower-level problem

$$\max_y \quad f(x,y) = y \quad (15a)$$

$$\text{s.t.} \quad -x + 2.5y \leq 3.75, \quad (15b)$$

$$x + 2.5y \geq 3.75, \quad (15c)$$

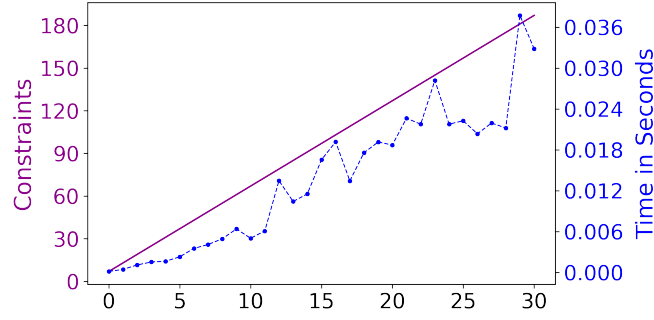
$$2.5x + y \leq 8.75, \quad (15d)$$

$$x, y \geq 0, \quad (15e)$$

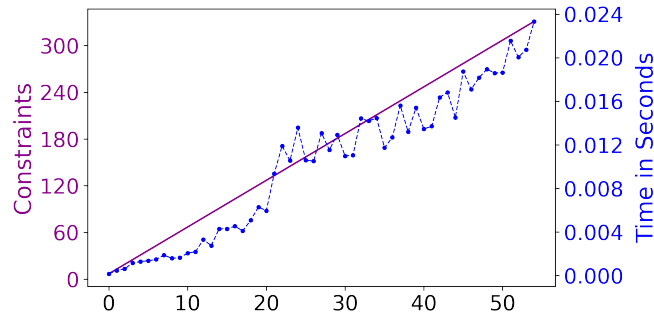
$$x, y \in \mathbb{R}. \quad (15f)$$

The optimal solution of Problem (14) is (0, 1.5). As explained in Section 5.1, we use the smallest and the largest x -values in the available data set to initialize $[\underline{x}, \bar{x}]$, which is required in Line 1 of Algorithm 1. For this instance, the unknown lower level is substituted by a neural network that has 2 hidden layers with 5 nodes each and is trained on a set of 30 points. The learning rate is approximately 0.074. Using the trained network's weights, we compute a Lipschitz constant of about 3.02. Finally, Algorithm 1 finds the approximate solution (0, 1.4999) in 1.21 seconds and 31 iterations for $\varepsilon = 10^{-5}$. Table 4 captures how different Lipschitz constants influence the computation time and the number of iterations of Algorithm 1 when applied to Problem (14). For all computations, we keep the tolerance $\varepsilon = 10^{-5}$. This table clearly shows that the algorithm performs better, the closer the used Lipschitz constant is to the true value of 2.5.

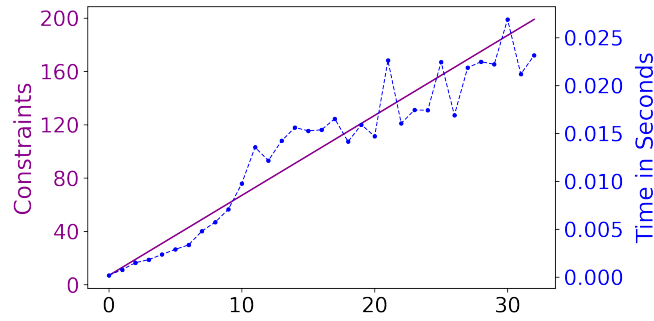
Moreover, constants smaller than 2.5 seemingly perform even better. For instance, using 0.5, Algorithm 1 finds a solution in only 10 iterations. Nevertheless, using values smaller than the Lipschitz constant of function $g_i(\cdot)$ can potentially lead to false infeasibilities reported by the algorithm.



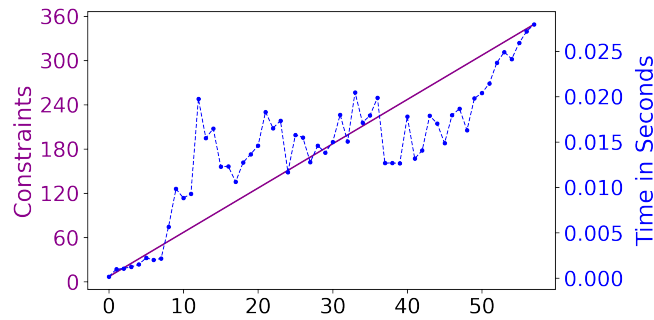
(A) Moore and Bard (1990)



(B) Bialas and Karwan (1984)



(C) Clark and Westerberg (1988)



(D) Liu and Hart (1994)

FIGURE 2. The computational effort of Algorithm 1 grows linearly in the number of constraints in each iteration. The x -axis shows the number of iterations required by Algorithm 1.

TABLE 4. Computation times in dependence of the Lipschitz constant. The second row corresponds to the Lipschitz constant obtained with LipSDP-Neuron, and the third one to the true Lipschitz constant of $\Psi(\cdot)$ for instance (14).

Lipschitz constant	Iterations	Computation time	Solution
5	65	3.67	(0, 1.4999)
3.02	31	1.08	(0, 1.4999)
2.5	33	1.06	(0, 1.5)
1.5	20	0.45	(0, 1.4999)
0.5	10	0.13	(0, 1.4999)
0.1	-	-	-

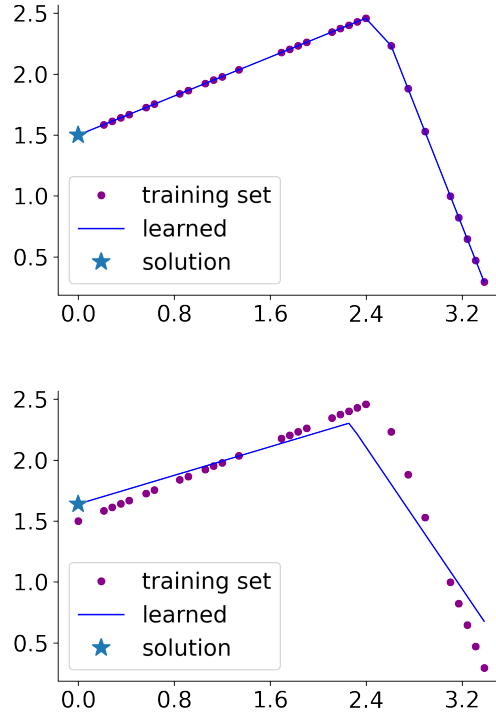


FIGURE 3. Two neural networks trained on the same set with different learning rates. Consequently, the learned function g (blue curves) is different and Algorithm 1 computes two different solutions (blue stars) for instance (14) using these functions.

5.2.2. *The Accuracy of Algorithm 1.* According to (11), an ε -feasible solution of an instance is ε -close to the graph of all functions $g_i(\cdot)$. This does not mean that it is also necessarily close to the true optimal response $\Psi_i(\cdot)$, $i \in [n_y]$. Consequently, the accuracy of the method depends on the ability of the neural network to approximate the true optimal responses as good as possible.

Figure 3 illustrates the importance of a good approximation. There you can see in the bottom plot that the ε -feasible solution, in this case (0.03, 1.67), obtained for instance (14), is ε -close to (0.03, $g(0.03)$) but not to (0.03, $\Psi(0.03)$).

TABLE 5. Lipschitz constants and solutions in dependence of the size of the training set for the instance in Bard (1984).

Training set size	Lipschitz constant	Solution
6	5.04	(6.5, 1.25)
12	6.64	(6.87, 1.43)
30	5.77	(7.07, 1.54)
60	6.35	(7.01, 1.5)
120	8.43	(7.17, 1.58)

In this context, it is also interesting that we observed that larger training sets can lead to better solutions. Table 5 shows how the computed Lipschitz constants and the solutions obtained for the instance in Bard (1984) change with increasing data sets. For this instance, the true solution is (7.2, 1.6). Due to the fact that the solution is located at an extreme point of \mathcal{F}_x , the accuracy of Algorithm 1 directly depends on the proxy for \mathcal{F}_x in this specific example.

6. CONCLUSION

In many practical situations, the leader of a bilevel optimization problem is not aware of an explicit formulation of the follower’s problem. To cope with this issue, we proposed a method that uses neural networks to learn the follower’s optimal reaction from past bilevel solutions. After training of the network, we compute Lipschitz constants for the learned functions and use a Lipschitz decomposition method to solve the reformulated, single-level problem with neural-network constraints.

This short paper should serve as a proof of concept for the ideas sketched above. However, many aspects can be improved. First of all, the assumption of a scalar leader’s decision seems rather strong. Very recently, a follow-up paper (Grübel et al. 2022) of Schmidt et al. (2019) appeared, in which a similar Lipschitz decomposition method is developed that can tackle the multi-dimensional case. This method can now be used to also consider bilevel optimization problems with an unknown follower problem and multiple decision variables of the leader. Second, we restricted ourselves in this paper to the case in which we have no coupling constraints. Fortunately, it is rather straightforward to use Algorithm 1 also for linear bilevel problems with coupling constraints as well. Even if the bilevel feasible set is disconnected due to the presence of coupling constraints, the follower’s optimal response function $\Psi_i(\cdot)$ is learned to be a piecewise linear function by the corresponding neural network. Due to the fact that solutions are found at vertices of the feasible set, at least one feasible point is always available to serve as solution, even if the line connecting two points is not bilevel feasible. On the other hand, this is getting more complicated if nonlinear bilevel problems are considered. Hence, the setting of nonlinear problems with coupling constraints is a topic of future research. Third, the overall idea should be reasonable also in the mixed-integer case, at least if no coupling constraints are present. Fourth, there has been some recent work (Bolusani et al. 2020) on the relation between multilevel mixed-integer linear optimization problems and multistage stochastic mixed-integer linear optimization problems with recourse. Hence, it might be possible to exploit these relations to carry over learning-based techniques for two-stage stochastic optimization to bilevel optimization. Fifth and finally, it would be very interesting to see an application of the concept discussed in this paper to a real-world situation, which is, however, out the of the scope of this short paper.

ACKNOWLEDGEMENTS

The authors thank the DFG for their support within RTG 2126 “Algorithmic Optimization”.

REFERENCES

- Bard, J. F. (1984). “Optimality conditions for the bilevel programming problem.” In: *Naval Research Logistics Quarterly* 31.1, pp. 13–26. DOI: [10.1002/nav.3800310104](https://doi.org/10.1002/nav.3800310104).
- Beck, Y., I. Ljubić, and M. Schmidt (2022). *A Survey on Bilevel Optimization Under Uncertainty*. Tech. rep. URL: http://www.optimization-online.org/DB_HTML/2022/06/8963.html.
- Besançon, M., M. F. Anjos, and L. Brotcorne (2019). *Near-optimal Robust Bilevel Optimization*. URL: <https://arxiv.org/pdf/1908.04040.pdf>.
- (2021). “Complexity of near-optimal robust versions of multilevel optimization problems.” In: *Optimization Letters* 15, pp. 2597–2610. DOI: [10.1007/s11590-021-01754-9](https://doi.org/10.1007/s11590-021-01754-9).
- Bialas, W. F. and M. H. Karwan (1984). “Two-level linear programming.” In: *Management Science* 30.8, pp. 1004–1020. URL: www.jstor.org/stable/2631591.
- Bolusani, S., S. Coniglio, T. K. Ralphs, and S. Tahernejad (2020). “A Unified Framework for Multistage Mixed Integer Linear Optimization.” In: *Bilevel Optimization: Advances and Next Challenges*. Ed. by S. Dempe and A. Zemkoho. Cham: Springer International Publishing, pp. 513–560. DOI: [10.1007/978-3-030-52119-6_18](https://doi.org/10.1007/978-3-030-52119-6_18).
- Borrero, J. S., O. A. Prokopyev, and D. Sauré (2022). “Learning in Sequential Bilevel Linear Programming.” In: *INFORMS Journal on Optimization*. DOI: [10.1287/ijoo.2021.0063](https://doi.org/10.1287/ijoo.2021.0063).
- Clark, P. and A. Westerberg (1988). “A note on the optimality conditions for the bilevel programming problem.” In: *Naval Research Logistics (NRL)* 35.5, pp. 413–418. DOI: [10.1002/1520-6750\(198810\)35:5<413::AID-NAV3220350505>3.0.CO;2-6](https://doi.org/10.1002/1520-6750(198810)35:5<413::AID-NAV3220350505>3.0.CO;2-6).
- Dempe, S. (2002). *Foundations of bilevel programming*. Springer Science & Business Media.
- (2020). “Bilevel Optimization: Theory, Algorithms, Applications and a Bibliography.” In: *Bilevel Optimization: Advances and Next Challenges*. Ed. by S. Dempe and A. Zemkoho. Springer International Publishing, pp. 581–672. DOI: [10.1007/978-3-030-52119-6_20](https://doi.org/10.1007/978-3-030-52119-6_20).
- Fazlyab, M., M. Morari, and G. J. Pappas (2022). “Safety Verification and Robustness Analysis of Neural Networks via Quadratic Constraints and Semidefinite Programming.” In: *IEEE Transactions on Automatic Control* 67.1, pp. 1–15. DOI: [10.1109/TAC.2020.3046193](https://doi.org/10.1109/TAC.2020.3046193).
- Fazlyab, M., A. Robey, H. Hassani, M. Morari, and G. Pappas (2019). “Efficient and Accurate Estimation of Lipschitz Constants for Deep Neural Networks.” In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc. URL: <https://proceedings.neurips.cc/paper/2019/file/95e1533eb1b20a9777749fb94fdb944-Paper.pdf>.
- Goodfellow, I., Y. Bengio, and A. Courville (2016). *Deep Learning*. MIT Press. URL: <https://mitpress.mit.edu/books/deep-learning>.
- Grübel, J., R. Krug, M. Schmidt, and W. Wollner (2022). *A Successive Linear Relaxation Method for MINLPs with Multivariate Lipschitz Continuous Nonlinearities with Applications to Bilevel Optimization and Gas Transport*. Tech. rep. URL: <http://arxiv.org/abs/2208.06444>.

- Hansen, P., B. Jaumard, and G. Savard (1992). “New branch-and-bound rules for linear bilevel programming.” In: *SIAM Journal on Scientific and Statistical Computing* 13.5, pp. 1194–1217. DOI: [10.1137/0913069](https://doi.org/10.1137/0913069).
- Haurie, A., G. Savard, and D. White (1990). “A note on: an efficient point algorithm for a linear two-stage optimization problem.” In: *Operations Research* 38.3, pp. 553–555. DOI: [10.1287/opre.38.3.553](https://doi.org/10.1287/opre.38.3.553).
- Khanduri, P., S. Zeng, M. Hong, H.-T. Wai, Z. Wang, and Z. Yang (2021). *A Near-Optimal Algorithm for Stochastic Bilevel Optimization via Double-Momentum*. URL: <https://arxiv.org/abs/2102.07367>.
- Kleinert, T., M. Labbé, I. Ljubić, and M. Schmidt (2021). “A Survey on Mixed-Integer Programming Techniques in Bilevel Optimization.” In: *EURO Journal of Computational Optimization* 9, p. 100007. DOI: [10.1016/j.ejco.2021.100007](https://doi.org/10.1016/j.ejco.2021.100007).
- Liu, Y.-H. and S. M. Hart (1994). “Characterizing an optimal solution to the linear bilevel programming problem.” In: *European Journal of Operational Research* 73.1, pp. 164–166. DOI: [10.1016/0377-2217\(94\)90155-4](https://doi.org/10.1016/0377-2217(94)90155-4).
- Lozano, L. and J. C. Smith (2017). “A Value-Function-Based Exact Approach for the Bilevel Mixed-Integer Programming Problem.” In: *Operations Research* 65.3, pp. 768–786. DOI: [10.1287/opre.2017.1589](https://doi.org/10.1287/opre.2017.1589).
- Moore, J. T. and J. F. Bard (1990). “The mixed integer linear bilevel programming problem.” In: *Operations Research* 38.5, pp. 911–921. URL: www.jstor.org/stable/171050.
- Pauli, P., A. Koch, J. Berberich, P. Kohler, and F. Allgöwer (2022). “Training Robust Neural Networks Using Lipschitz Bounds.” In: *IEEE Control Systems Letters* 6, pp. 121–126. DOI: [10.1109/LCSYS.2021.3050444](https://doi.org/10.1109/LCSYS.2021.3050444).
- Rios, L. M. and N. V. Sahinidis (2013). “Derivative-free optimization: a review of algorithms and comparison of software implementations.” In: *Journal of Global Optimization* 56.3, pp. 1247–1293. DOI: [10.1007/s10898-012-9951-y](https://doi.org/10.1007/s10898-012-9951-y).
- Ruder, S. (2016). *An overview of gradient descent optimization algorithms*. URL: <https://arxiv.org/abs/1609.04747>.
- Schmidt, M., M. Sirvent, and W. Wollner (2019). “A Decomposition Method for MINLPs with Lipschitz Continuous Nonlinearities.” In: *Mathematical Programming* 178.1, pp. 449–483. DOI: [10.1007/s10107-018-1309-x](https://doi.org/10.1007/s10107-018-1309-x).
- (2021). “The Cost of Not Knowing Enough: Mixed-Integer Optimization with Implicit Lipschitz Nonlinearities.” In: *Optimization Letters*. DOI: [10.1007/s11590-021-01827-9](https://doi.org/10.1007/s11590-021-01827-9). Forthcoming.
- Vlah, D., K. Šepetanc, and H. Pandžić (2022). *Solving Bilevel Optimal Bidding Problems Using Deep Convolutional Neural Networks*. DOI: [10.48550/ARXIV.2207.05825](https://doi.org/10.48550/ARXIV.2207.05825).

(I. Molan, M. Schmidt) TRIER UNIVERSITY, DEPARTMENT OF MATHEMATICS, UNIVERSITÄTSRING
15, 54296 TRIER, GERMANY
Email address: molan@uni-trier.de
Email address: martin.schmidt@uni-trier.de