

LEVENBERG–MARQUARDT METHOD AND PARTIAL EXACT PENALTY PARAMETER SELECTION IN BILEVEL OPTIMIZATION

Andrey Tin[†] and Alain B. Zemkoho[‡]

*Centre for Operational Research, Management Sciences and Information Systems (CORMSIS)
and School of Mathematical Sciences, University of Southampton, SO17 1BJ Southampton, UK*

January 26, 2021

ABSTRACT

We consider the optimistic bilevel optimization problem, known to have a wide range of applications in engineering, that we transform into a single-level optimization problem by means of the lower-level optimal value function reformulation. Subsequently, based on the partial calmness concept, we build an equation system, which is parameterized by the corresponding partial exact penalization parameter. We then design and analyze a Levenberg-Marquardt method to solve this parametric system of equations. Considering the fact that the selection of the partial exact penalization parameter is a critical issue when numerically solving a bilevel optimization problem, we conduct a careful experimental study to this effect, in the context the Levenberg-Marquardt method, while using the Bilevel Optimization LIBrary (BOLIB) series of test problems.

1 INTRODUCTION

We consider the optimistic bilevel optimization problem

$$\min_{x,y} F(x,y) \text{ s.t. } G(x,y) \leq 0, y \in S(x) := \arg \min_y \{f(x,y) : g(x,y) \leq 0\}, \quad (1.1)$$

where $F : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$, $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$, $G : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^q$, and $g : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^p$. As usual, we refer to F (resp. f) as the upper-level (resp. lower-level) objective function and G (resp. g) stands for the upper-level (resp. lower-level) constraint function. Solving problem (1.1) is very difficult because of the implicit nature of the lower-level optimal solution set-valued mapping $S : \mathbb{R}^n \rightrightarrows \mathbb{R}^m$ that partly describes the feasible set of problem (1.1). Problem (1.1) has a wide range of applications in engineering; to have a flavour of this, see, e.g., [7] and references therein.

One of the most common single-level transformation of problem (1.1) is the Karush-Kuhn-Tucker (KKT) reformulation, which consists of replacing inclusion $y \in S(x)$ with the equivalent KKT conditions of the lower-level problem, under appropriate assumptions; see, e.g., [1, 8, 5, 15, 24]. As shown in [5], the first drawback of the KKT reformulation is that it is not necessarily equivalent to the original problem (1.1). Secondly, from the numerical perspective, the KKT reformulation involves derivatives from the lower-level problem that can require the calculation of second (resp.

[†] e-mail: a.tin@soton.ac.uk.

[‡] e-mail: a.b.zemkoho@soton.ac.uk.

AT was supported by a University of Southampton Vice Chancellor Scholarship, while ABZ was supported by the EPSRC grant EP/V049038/1 and the Alan Turing Institute under the EPSRC grant EP/N510129/1.

third) order derivatives for first (resp. second) order optimization-type methods. Thirdly, it is shown in the new paper [31] that the lower-level value function (LLVF) reformulation

$$\min_{x,y} F(x,y) \text{ s.t. } G(x,y) \leq 0, \quad g(x,y) \leq 0, \quad f(x,y) \leq \varphi(x), \quad (1.2)$$

where the LLVF φ is defined by

$$\varphi(x) := \min \{f(x,y) \mid g(x,y) \leq 0\}, \quad (1.3)$$

can lead to a better numerical performance than the KKT reformulation for certain problem classes; in particular for the Bilevel Optimization LIBrary (BOLIB) examples [33]. It is also important note that problem (1.2)–(1.3) is equivalent to the original problem (1.1) without any assumption.

There is a number of recent studies on solution methods for bilevel programs, based on the LLVF reformulation. For example, [20, 22, 27] develop global optimization techniques for (1.2)–(1.3). [18, 28] propose algorithms computing stationary points for (1.2)–(1.3) in the special case where G (resp. g) is independent from y (resp. x). But, the value function in the latter work is approximated by an entropy function, which is difficult to compute for problems with a big number of lower-level variables, as the corresponding approximation relies on integral calculations.

More recently, [12, 31] have proposed Newton-type methods for problem (1.2)–(1.3) based on a special transformation that enables the optimality conditions of this problem to be squared. Naturally, detailed optimality conditions for (1.2)–(1.3) are non-square and are directly dealt with in [14] with Gauss-Newton-type methods. But, these methods require inverse calculations at each iteration for matrices which are not necessarily nonsingular. Hence, to expand the applicability of this method to a wider class of problems, we consider and study the convergence of a regularized version in this paper, commonly known as the Levenberg-Marquardt method.

Recall that the standard approach to derive optimality conditions for problem (1.2)–(1.3), necessary to build this Levenberg-Marquardt method, is based on the concept of partial calmness. The consequence of this is that all the methods studied in the papers [14, 12, 31, 18, 28] are based on a partial exact penalization parameter, which needs a special care to be selected. However, none of these papers pays a special attention on how to select this parameter. One of the goals of the current paper is to conduct a numerical study on the best way to select this parameter, based on the BOLIB test set [33], and in the context of the Levenberg-Marquardt method to be introduced in the next section. This study will certainly enlightening the process of selection the partial exact penalization parameter in a broader class of methods to solve optimistic bilevel optimization problems.

Another difficulty working with the LLVF reformulation is that φ is typically non-differentiable function. This will be handled in this paper by using upper estimates of the subdifferential of the function (see, e.g., [9, 6, 8, 30]) that will lead to a relatively simple system of optimality conditions, which depends on the aforementioned partial exact penalization parameter.

The remainder of the paper proceeds as follows. In the next section, we start by transforming the necessary optimality conditions into a system of equations. To do so, we substitute the corresponding complementarity conditions by the well-known Fischer-Burmeister function [11] that we subsequently perturb to build a smooth system of equations. In Section 3, we focus our attention on partial exact penalization, including a detailed discussion on how to chose the penalty parameter, as well as the related ill-behaviour. This analysis plays an important role in the performance of the method, given that a particular attention is paid on ways to avoid the aforementioned ill-behaviour. Furthermore, we will present results on the experimental order of convergence of the method and line search stepsize parameter at the last iteration. The results in Section 4 are compared with known solutions of the problems to check the performance of the method under different partial exact penalization parameter scenarios.

2 THE ALGORITHM AND CONVERGENCE ANALYSIS

We start this section with some definitions necessary to state optimality conditions for the bilevel program (1.2)–(1.3), which will be used to construct the algorithm. The lower-level problem is *fully*

convex if the functions f and g_i , $i = 1, \dots, p$ are convex in (x, y) . A point (\bar{x}, \bar{y}) is said to be *lower-level regular* if there exists a vector $d \in \mathbb{R}^m$ such that we have

$$\nabla_y g_i(\bar{x}, \bar{y})^\top d < 0 \quad \text{for } i \in I_g(\bar{x}, \bar{y}) := \{i : g_i(\bar{x}, \bar{y}) = 0\}. \quad (2.1)$$

Obviously, this is the *Mangasarian-Fromovitz constraint qualification* (MFCQ) for the lower-level problem in (1.1). Similarly, a point (\bar{x}, \bar{y}) is *upper-level regular* if there exists $d \in \mathbb{R}^n \times \mathbb{R}^m$ such that

$$\begin{aligned} \nabla g_i(\bar{x}, \bar{y})^\top d &< 0 \quad \text{for } j \in I_g(\bar{x}, \bar{y}), \\ \nabla G_j(\bar{x}, \bar{y})^\top d &< 0 \quad \text{for } j \in I_G(\bar{x}, \bar{y}) := \{j : G_j(\bar{x}, \bar{y}) = 0\}. \end{aligned} \quad (2.2)$$

Finally, to write the necessary optimality conditions for problem (1.2)–(1.3), it is standard to use the following partial calmness concept [30]:

Definition 2.1. *Let (\bar{x}, \bar{y}) be a local optimal solution of problem (1.2)–(1.3). The problem is partially calm at (\bar{x}, \bar{y}) if there exists $\lambda > 0$ and a neighbourhood U of $(\bar{x}, \bar{y}, 0)$ such that*

$$F(x, y) - F(\bar{x}, \bar{y}) + \lambda|u| \geq 0, \quad \forall (x, y, u) \in U : G(x, y) \leq 0, \quad g(x, y) \leq 0, \quad f(x, y) - \varphi(x) - u = 0.$$

The following relationship shows that the partial calmness concept enables the penalization of the optimal value function constraint $f(x, y) - \varphi(x) \leq 0$, to generate a tractable optimization problem, as it is well-known that standard constraint qualifications (including the MFCQ, for example) fail for problem (1.2)–(1.3); see [9, 30].

Theorem 2.2 ([30]). *Let (\bar{x}, \bar{y}) be a local minimizer of (1.2)–(1.3). Then, this problem is partially calm at (\bar{x}, \bar{y}) if and only if there is some $\bar{\lambda} > 0$ such that for any $\lambda \geq \bar{\lambda}$, the point (\bar{x}, \bar{y}) is also a local minimizer of*

$$\min_{x, y} F(x, y) + \lambda(f(x, y) - \varphi(x)) \quad \text{s.t. } G(x, y) \leq 0, \quad g(x, y) \leq 0. \quad (2.3)$$

Problem (2.3) is known as the *partial exact penalization* problem of (1.2)–(1.3), as only the optimal value constraint is penalized. Next, we state the necessary optimality conditions for problem (1.2)–(1.3) based on the aforementioned penalization, while using the upper- and lower-level regularity conditions and an estimate of the subdifferential of φ (1.3); see, e.g., [9, 6, 8, 30], for the details.

Theorem 2.3. *Let (\bar{x}, \bar{y}) be a local optimal solution of problem (1.2), where all involved functions are assumed to be continuously differentiable, φ is finite around \bar{x} , and the lower-level problem is fully convex. Furthermore, suppose that the problem is partially calm at (\bar{x}, \bar{y}) , and the lower- and upper-level regularity conditions are both satisfied at (\bar{x}, \bar{y}) . Then, there exist $\lambda > 0$, as well as $u, w \in \mathbb{R}^p$ and $v \in \mathbb{R}^q$ such that*

$$\nabla F(\bar{x}, \bar{y}) + \nabla g(\bar{x}, \bar{y})^\top (u - \lambda w) + \nabla G(\bar{x}, \bar{y})^\top v = 0, \quad (2.4)$$

$$\nabla_y f(\bar{x}, \bar{y}) + \nabla_y g(\bar{x}, \bar{y})^\top w = 0, \quad (2.5)$$

$$u \geq 0, \quad g(\bar{x}, \bar{y}) \leq 0, \quad u^\top g(\bar{x}, \bar{y}) = 0, \quad (2.6)$$

$$v \geq 0, \quad G(\bar{x}, \bar{y}) \leq 0, \quad v^\top G(\bar{x}, \bar{y}) = 0, \quad (2.7)$$

$$w \geq 0, \quad g(\bar{x}, \bar{y}) \leq 0, \quad w^\top g(\bar{x}, \bar{y}) = 0. \quad (2.8)$$

In this result, partial calmness and full convexity are essential and fundamentally related to the nature of the bilevel optimization. Hence, it is important to highlight a few classes of problems satisfying these assumptions. Partial calmness has been the main tool to derive optimality conditions for (1.1) from the perspective of the optimal value function; see, e.g., [6, 8, 9, 30]. It automatically holds if G is independent from y and the lower-level problem is defined by

$$f(x, y) := c^\top y \quad \text{and} \quad g(x, y) := A(x) + By, \quad (2.9)$$

where $A : \mathbb{R}^n \rightarrow \mathbb{R}^p$, $c \in \mathbb{R}^m$, and $B \in \mathbb{R}^{p \times m}$. More generally, various sufficient conditions ensuring that partial calmness holds have been studied in the literature; see [30] for the seminal work on the subject. More recently, the paper [19] has revisited the condition, proposed a fresh perspective, and established new sufficient conditions for it to be satisfied.

As for full convexity, it will automatically hold for the class of problem defined in (2.9), provided that each component of the function A is convex. Another class of nonlinear fully convex lower-level problem is given in [17]. Note however that when it is not possible to guarantee that this assumption is satisfied, there are at least two alternative scenarios to obtain the same optimality conditions as in Theorem 2.3. The first is to replace the full convexity assumption by the *inner semicontinuity* of the optimal solution set-valued mapping S (1.1). Secondly, note that a much weaker qualification condition known as *inner semicompactness* can also be used here. However, under the latter assumption, it will additionally be required to have $S(\bar{x}) = \{\bar{y}\}$ in order to get the optimality conditions in (2.4)–(2.8). The concept of inner semicontinuity (resp. semicompactness) of S is closely related to the lower semicontinuity (resp. upper semicontinuity) of set-valued mappings; for more details on these notions and their ramifications on bilevel programs, see, e.g., [6, 8, 9].

It is important to mention that various other necessary optimality conditions, different from the above ones, can be obtained, depending on the assumptions made. Details of different stationarity concepts for (1.2) can be found in the latter references, as well as in [32].

To reformulate the complementarity conditions (2.6)–(2.8) into a system of equations, we use the well-known *Fischer-Burmeister* function [11] $\phi(a, b) := \sqrt{a^2 + b^2} - a - b$, which ensures that

$$\phi(a, b) = 0 \iff a \geq 0, \quad b \geq 0, \quad ab = 0.$$

This leads to the reformulation of the optimality conditions (2.4)–(2.8) into the system of equations:

$$\Upsilon^\lambda(z) := \begin{pmatrix} \nabla_x F(x, y) + \nabla_x g(x, y)^T(u - \lambda w) + \nabla_x G(x, y)^T v \\ \nabla_y F(x, y) + \nabla_y g(x, y)^T(u - \lambda w) + \nabla_y G(x, y)^T v \\ \nabla_y f(x, y) + \nabla_y g(x, y)^T w \\ \sqrt{u^2 + g(x, y)^2} - u + g(x, y) \\ \sqrt{v^2 + G(x, y)^2} - v + G(x, y) \\ \sqrt{w^2 + g(x, y)^2} - w + g(x, y) \end{pmatrix} = 0, \quad (2.10)$$

where we have $z := (x, y, u, v, w)$ and

$$\sqrt{u^2 + g(x, y)^2} - u + g(x, y) := \begin{pmatrix} \sqrt{u_1^2 + g_1(x, y)^2} - u_1 + g_1(x, y) \\ \vdots \\ \sqrt{u_p^2 + g_p(x, y)^2} - u_p + g_p(x, y) \end{pmatrix}. \quad (2.11)$$

$\sqrt{v^2 + G(x, y)^2} - v + G(x, y)$ and $\sqrt{w^2 + g(x, y)^2} - w + g(x, y)$ are defined as in (2.11). The superscript λ is used to emphasize the fact that this number is a parameter and not a variable for equation (2.10). One can easily check that this system is made of $n + 2m + p + q + p$ real-valued equations and $n + m + p + q + p$ variables. Clearly, this means that (2.10) is an *overdetermined* system and the Jacobian of $\Upsilon^\lambda(z)$, when it exists, is a non-square matrix.

To focus our attention on the main ideas of this paper, we smoothen the function Υ^λ (2.10) using the *smoothed Fischer-Burmeister function* (see [16]) defined by

$$\phi_j^\mu(x, y, u) := \sqrt{u_j^2 + g_j(x, y)^2 + 2\mu} - u_j + g_j(x, y), \quad j = 1, \dots, p, \quad (2.12)$$

where the perturbation $\mu > 0$ helps to guaranty its differentiability at points (x, y, u) , where $u_j = g_j(x, y) = 0$ for $j = 1, \dots, p$. It is well-known (see latter reference) that for $j = 1, \dots, p$,

$$\phi_j^\mu(x, y, u) = 0 \iff [u_j > 0, \quad -g_j(x, y) > 0, \quad -u_j g_j(x, y) = \mu]. \quad (2.13)$$

The smoothed version of system (2.10) then becomes

$$\Upsilon_\mu^\lambda(z) = \begin{pmatrix} \nabla_x F(x, y) + \nabla_x g(x, y)^T(u - \lambda w) + \nabla_x G(x, y)^T v \\ \nabla_y F(x, y) + \nabla_y g(x, y)^T(u - \lambda w) + \nabla_y G(x, y)^T v \\ \nabla_y f(x, y) + \nabla_y g(x, y)^T w \\ \sqrt{u^2 + g(x, y)^2 + 2\mu} - u + g(x, y) \\ \sqrt{v^2 + G(x, y)^2 + 2\mu} - v + G(x, y) \\ \sqrt{w^2 + g(x, y)^2 + 2\mu} - w + g(x, y) \end{pmatrix} = 0, \quad (2.14)$$

following the convention in (2.11), where μ is a vector of appropriate dimensions with sufficiently small positive elements. Under the assumption that all the functions involved in problem (1.1) are continuously differentiable, Υ_μ^λ is also a continuously differentiable function for any $\lambda > 0$ and $\mu > 0$. We can easily check that for a fixed value of $\lambda > 0$,

$$\|\Upsilon_\mu^\lambda(z) - \Upsilon^\lambda(z)\| \rightarrow 0 \text{ as } \mu \downarrow 0. \quad (2.15)$$

Hence, based on this scheme, our aim is to consider a sequence $\{\mu_k\}$ decreasing to 0 such that equation (2.10) is approximately solved while leading to

$$\Upsilon_{\mu^k}^\lambda(z^k) \rightarrow 0 \text{ as } k \uparrow \infty$$

for a fixed value of $\lambda > 0$. In order to proceed, let us define the least squares problem

$$\min \Phi_\mu^\lambda(z) := \frac{1}{2} \|\Upsilon^\lambda(z)\|^2. \quad (2.16)$$

Before we introduce the smoothed Levenberg-Marquardt method that will be one of the main focus points of this paper, note that for fixed $\lambda > 0$ and $\mu > 0$,

$$\nabla \Upsilon_\mu^\lambda(z) = \begin{bmatrix} \nabla^2 L^\lambda(z) & \nabla g(x, y)^T & \nabla G(x, y)^T & -\lambda \nabla g(x, y)^T \\ \nabla(\nabla_y \mathcal{L}(z)) & \mathbf{O} & \mathbf{O} & \nabla_y g(x, y)^T \\ \mathcal{J}^\mu \nabla g(x, y) & \Gamma^\mu & \mathbf{O} & \mathbf{O} \\ \mathcal{A}^\mu \nabla G(x, y) & \mathbf{O} & \mathcal{B}^\mu & \mathbf{O} \\ \Theta^\mu \nabla g(x, y) & \mathbf{O} & \mathbf{O} & \mathcal{K}^\mu \end{bmatrix} \quad (2.17)$$

with the pair $(\mathcal{J}^\mu, \Gamma^\mu)$ defined by $\mathcal{J}^\mu := \text{diag}\{\tau_1^\mu, \dots, \tau_p^\mu\}$ and $\Gamma^\mu := \text{diag}\{\gamma_1^\mu, \dots, \gamma_p^\mu\}$, where

$$\tau_j^\mu := \frac{g_j(x, y)}{\sqrt{u_j^2 + g_j(x, y)^2 + 2\mu}} + 1 \text{ and } \gamma_j^\mu := \frac{u_j}{\sqrt{u_j^2 + g_j(x, y)^2 + 2\mu}} - 1, \quad j = 1, \dots, p. \quad (2.18)$$

The pairs $(\mathcal{A}^\mu, \mathcal{B}^\mu)$ and $(\Theta^\mu, \mathcal{K}^\mu)$ are defined in a similar way, based on $(v_j, G_j(x, y))$, $j = 1, \dots, q$ and $(w_j, g_j(x, y))$, $j = 1, \dots, p$, respectively.

We now move on to present some definitions that will help us state the algorithm with *line search*. It is well-known that line search helps to choose the optimal step length to avoid over-going an optimal solution in the direction d^k and also to globalize the convergence of the method, i.e., have more flexibility on the starting point z^0 . The optimal step length γ_k can be calculated through minimizing $\Phi^\lambda(z^k + \gamma_k d^k)$, with respect to γ_k , such that

$$\Phi^\lambda(z^k + \gamma_k d^k) \leq \Phi^\lambda(z^k) + \sigma \gamma_k \nabla \Phi^\lambda(z^k)^T d^k \text{ for } 0 < \sigma < 1.$$

That is, we are looking for $\gamma_k = \arg \min_{\gamma \in \mathbb{R}} \|\Upsilon^\lambda(z^k + \gamma d^k)\|^2$. To implement line search, it is standard to use *Armijo condition* that guarantees a decrease at the next iterate.

Definition 2.4. Fixing d and z , consider the function $\phi_\lambda(\gamma) := \Phi^\lambda(z + \gamma d)$. Then, the Armijo condition will be said to hold if $\phi_\lambda(\gamma) \leq \phi(0) + \gamma \sigma \phi'_\lambda(\gamma)$ for some $0 < \sigma < 1$.

The practical implementation of the Armijo condition is based on backtracking.

Definition 2.5. Let $\rho \in (0, 1)$ and $\bar{\gamma} > 0$. Backtracking is the process of checking over a sequence $\bar{\gamma}, \rho \bar{\gamma}, \rho^2 \bar{\gamma}, \dots$, until a number γ is found satisfying the Armijo condition.

Line search is widely used in continuous optimization; see, e.g., [21] for more details. For the implementation in this paper, we start with stepsize $\gamma_0 := 1$; then, if the condition

$$\|\Upsilon^\lambda(z^k + \gamma_k d^k)\|^2 < \|\Upsilon^\lambda(z^k)\|^2 + \sigma \gamma_k \nabla \Upsilon_\mu^\lambda(z^k)^T \Upsilon^\lambda(z^k) d^k,$$

is not satisfied, we set $\gamma_k = \gamma_k/2$ and check again until the condition above is satisfied. More generally, the algorithm proceeds as follows:

Algorithm 2.6. *Smoothed Levenberg-Marquardt Method for Bilevel Optimization*

Step 0: Choose $(\lambda, \mu, K, \epsilon, \alpha_0) \in (\mathbb{R}_+^*)^5$, $(\rho, \sigma, \gamma_0) \in (0, 1)^3$, $z^0 := (x^0, y^0, u^0, v^0, w^0)$, and set $k := 0$.

Step 1: If $\|\Upsilon_\mu^\lambda(z^k)\| < \epsilon$ or $k \geq K$, then stop.

Step 2: Calculate the Jacobian $\nabla\Upsilon_\mu^\lambda(z^k)$ and subsequently find a vector d^k satisfying

$$\left(\nabla\Upsilon_\mu^\lambda(z^k)^\top \nabla\Upsilon_\mu^\lambda(z^k) + \alpha_k I\right) d^k = -\nabla\Upsilon_\mu^\lambda(z^k)^\top \Upsilon_\mu^\lambda(z^k), \quad (2.19)$$

where I denotes the identity matrix of appropriate size.

Step 3: While $\|\Upsilon_\mu^\lambda(z^k + \gamma_k d^k)\|^2 \geq \|\Upsilon_\mu^\lambda(z^k)\|^2 + \sigma \gamma_k \nabla\Upsilon_\mu^\lambda(z^k)^\top \Upsilon_\mu^\lambda(z^k) d^k$, do $\gamma_k \leftarrow \rho \gamma_k$ end.

Step 4: Set $z^{k+1} := z^k + \gamma_k d^k$, $k := k + 1$, and go to Step 1.

Note that in Step 0, $\mathbb{R}_+^* := (0, \infty)$. Before we move on to focus our attention on the practical implementation details of this algorithm, we present its convergence result, which is based on the following selection of the Levenberg–Marquardt (LM) parameter α_k :

$$\alpha_k = \|\Upsilon_\mu^\lambda(z^k)\|^\eta \text{ for any choice of } \eta \in [1, 2]. \quad (2.20)$$

Theorem 2.7 ([10]). *Consider Algorithm 2.6 with fixed values for the parameters $\lambda > 0$ and $\mu > 0$ and let α_k be defined as in (2.20). Then, the sequence $\{z^k\}$ generated by the algorithm converges quadratically to \bar{z} satisfying $\Upsilon_\mu^\lambda(\bar{z}) = 0$, under the following assumptions:*

1. $\Upsilon_\mu^\lambda : \mathbb{R}^N \rightarrow \mathbb{R}^{N+m}$ is continuously differentiable and $\nabla\Upsilon_\mu^\lambda : \mathbb{R}^N \rightarrow \mathbb{R}^{(N+m) \times (N)}$ is locally Lipschitz continuous in a neighbourhood of \bar{z} .
2. There exists some $C > 0$ and $\delta > 0$ such that

$$C \text{dist}(z, Z_\mu^\lambda) \leq \|\Upsilon_\mu^\lambda(z)\| \text{ for all } z \in \mathcal{B}(\bar{z}, \delta),$$

where dist denotes the distance function and Z_μ^λ corresponds to the solution set of equation (2.14).

For fixed values of $\lambda > 0$ and $\mu > 0$, assumption 1 in this theorem is automatically satisfied if all the functions involved in problem (1.1) are twice continuously differentiable. According to [29], assumption 2 of Theorem 2.7 is fulfilled if the matrix $\nabla\Upsilon_\mu^\lambda$ has a full column rank. Various conditions guarantying that $\nabla\Upsilon_\mu^\lambda$ has a full rank have been developed in [14]. Below, we present an example of bilevel program satisfying the first and second assumptions of Theorem 2.7.

Example 2.8. Consider the following instance of problem (1.1) from the BOLIB library [33, LamprielloSagratelli2017Ex33]:

$$\begin{aligned} F(x, y) &:= x^2 + (y_1 + y_2)^2, \\ G(x, y) &:= -x + 0.5, \\ f(x, y) &:= y_1, \\ g(x, y) &:= \begin{pmatrix} -x - y_1 - y_2 + 1 \\ -y \end{pmatrix}. \end{aligned}$$

Obviously, assumption 1 holds. According to [14], for this example, the columns of $\nabla\Upsilon_\mu^\lambda$ are linearly independent at the solution point

$$\bar{z} := (\bar{x}, \bar{y}_1, \bar{y}_2, \bar{u}_1, \bar{u}_2, \bar{u}_3, \bar{v}, \bar{w}_1, \bar{w}_2, \bar{w}_3) = (0.5, 0, 0.5, 1, \lambda, 0, 0, 0, 1, 0)$$

with the parameters chosen as $\mu = 2 \times 10^{-2}$ and $\lambda = 10^{-2}$. \square

On the selection of the LM parameter α_k , we conducted a preliminary analysis based on the BOLIB library test set [33]. It was observed that for almost all the corresponding examples, the choice $\alpha_k := \|\Upsilon_\mu^\lambda(z^k)\|^\eta$ with $\eta \in (1, 2]$ leads to a very poor performance of Algorithm 2.6. The typical behaviour of the algorithm for $\eta \in (1, 2]$ is shown in the following example.

Example 2.9. Consider the following scenario of problem (1.1) from [33, AllendeStill2013]:

$$\begin{aligned} F(x, y) &:= x_1^2 - 2x_1 + x_2^2 - 2x_2 + y_1^2 + y_2^2, \\ G(x, y) &:= \begin{pmatrix} -x \\ -y \\ x_1 - 2 \end{pmatrix}, \\ f(x, y) &:= y_1^2 - 2x_1y_1 + y_2^2 - 2x_2y_2, \\ g(x, y) &:= \begin{pmatrix} (y_1 - 1)^2 - 0.25 \\ (y_2 - 1)^2 - 0.25 \end{pmatrix}. \end{aligned}$$

Figure 1 shows the progression of $\|\Upsilon^\lambda(z^k)\|$ generated from Algorithm 2.6 with α_k selected as in (2.20) while setting $\eta = 1$ and $\eta = 2$, respectively. Clearly, after about 100 iterations $\|\Upsilon^\lambda(z^k)\|$ blows up relatively quickly when $\eta = 2$, while it falls and stabilizes within a certain tolerance for $\eta = 1$.

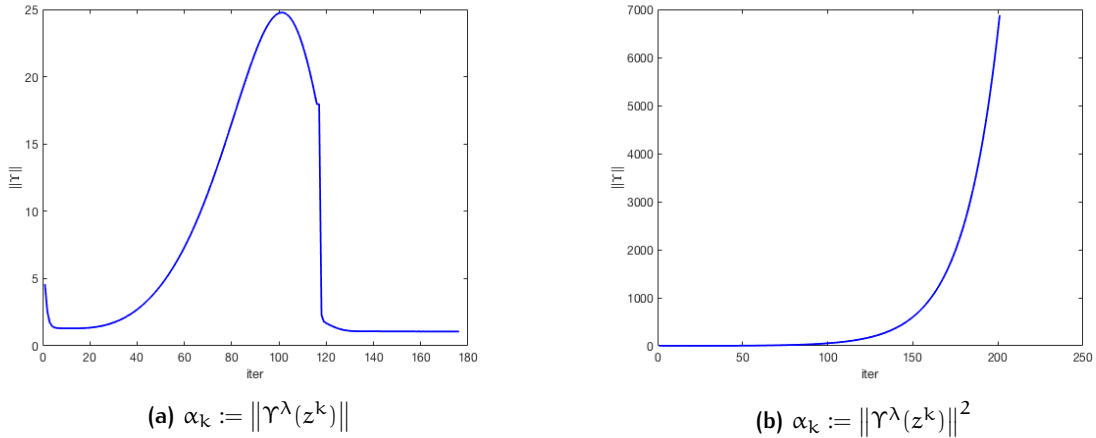


Figure 1: Typical behaviour of Algorithm 2.6 for two scenarios of the LM parameter

It is worth noting that the scale on the y-axis of Figure 1(b) is quite large. Hence, it might not be apparent that solutions at the early iterations of the algorithm are much better for $\eta = 1$ compared to the ones obtained in the case where $\eta = 2$. \square

Note that for almost all of the examples in the BOLIB test set [33], we have a behaviour of Algorithm 2.6 similar to that of Figure 1(b) when $\alpha_k = \|\Upsilon^\lambda(z^k)\|^\eta$ for many different choices of $\eta \in (1, 2]$. It is important to mention that the behaviour of the algorithm is not surprising, as it is well-known in the literature that with the sequence $\alpha_k := \|\Upsilon_\mu^\lambda(z^k)\|^2$, for example, the Levenberg–Marquardt method often faces some issues. Namely, when the sequence z^k is close to the solution set, $\alpha_k := \|\Upsilon_\mu^\lambda(z^k)\|^2$ could become smaller than the machine precision and hence lose its role as a result of this. On the other hand, when the sequence is far away from the solution set, $\|\Upsilon_\mu^\lambda(z^k)\|^2$ may be very large; making a movement towards the solution set to be very slow. Hence, from now on, we use $\alpha_k := \|\Upsilon^\lambda(z^k)\|$, with $k = 0, 1, \dots$, for all the analysis of Algorithm 2.6 conducted in this paper. Note however that there are various other approaches to select the LM parameter α_k in the literature; see, e.g., [2, 10, 29].

To conclude this section, it is important to recall that from the perspective of bilevel optimization, the partial exact penalization parameter λ is a key element of Algorithm 2.6, as it originates from the penalization of the value function constraint $f(x, y) \leq \varphi(x)$. Additionally, unlike the other parameters involved in the algorithm, which have benefited from many years of research, as it will be clear in Section 4, it remains unknown what is the best way to select λ while solving problem (1.1) via the value function reformulation (1.2)–(1.3). Hence, the focus of the remaining parts of this paper will be on the selection of λ and assessing its impact on the efficiency of Algorithm 2.6.

3 PARTIAL EXACT PENALTY PARAMETER SELECTION

The aim of this section is to explore the best way to select the penalization parameter λ . Based on Theorem 2.2, we should be getting the solution for some threshold value $\bar{\lambda}$ of the penalty parameter and the algorithm should be returning the value of the solution for any $\lambda \geq \bar{\lambda}$. Hence, increasing λ at every iteration seems to be the ideal approach to follow this logic to obtain and retain the solution. Hence, to proceed, we use the increasing sequence $\lambda_k := 0.5 * (1.05)^k$, where k is the number of iterations of the algorithm. The main reason of this choice is that the final value of λ need not to be too small to recover solutions and not too large to avoid ill-conditioning. It was observed that going too aggressive with the sequence (e.g., with $\lambda_k := 2^k$) forces the algorithm to diverge. Also, it was observed that for fixed and small values of λ (i.e., $\lambda < 1$), Algorithm 2.6 performs well for many examples. This justifies choosing the starting value for varying parameter less than 1. Overall, the aim here is vary λ as mentioned above and assess what could be the potential best ranges for selection of the parameter based on our test set from BOLIB [33].

3.1 How far can we go with the value of λ ?

We start here by acknowledging that it is very difficult to check that partial calmness (cf. Definition 2.1) holds in practice. Nevertheless, we would like to consider Theorem 2.2 as the base of our experimental analysis, and ask ourselves how large does λ need to be for Algorithm 2.6 to converge. Intuitively, one would think that taking λ as whatever large number should be fine. However, this is usually not the case in practice. One of the main reasons for this is that for too large values of λ , Algorithm 2.6 does not behave well. In particular, if we run Algorithm 2.6 with varying λ for a very large number of iterations, the value of the Error blows up at some point and the values $\|\gamma_\mu^\lambda(z^k)\|$ stop decreasing. Recall that ill-conditioning (and possibly the *ill-behaviour*) refers to the fact that one eigenvalue of the Hessian involved in $\nabla \gamma_\mu^\lambda(z^k)$ being much larger than the other eigenvalues, which affects the curvature in the negative way for gradient methods [25]. To analyze this ill-behaviour in this section, we are going to run our algorithm for 1,000 iterations with no stopping criteria and let λ vary indefinitely. We will subsequently look at which iteration the algorithm blows up and record the value of λ there. To proceed, we denote by λ_{ill} the first value of λ for which the ill-behaviour is observed for each example. The table below presents the number of BOLIB examples that are approximately within certain intervals of λ_{ill} .

Table 1: Approximate intervals for the values of λ when the ill-behaviour starts

λ_{ill}	$\lambda_{\text{ill}} < 10^7$	$10^7 < \lambda_{\text{ill}} < 10^9$	$10^9 < \lambda_{\text{ill}} < 10^{11}$	$10^{11} < \lambda_{\text{ill}} < 10^{20}$	Not observed
Examples	1	6	72	11	34

On average, the ill-behaviour seems to typically occur after about 500 iterations (where $\lambda_{\text{ill}} \approx 10^{10}$), as it can be seen in the table above. For 34 problems ill behaviour was not observed under the scope of 1000 iterations. We also see that for most of the problems (72/124), the ill-behaviour happens for $10^9 < \lambda < 10^{11}$. We further observe that algorithm has shown to behave well for the values of penalty parameter $\lambda < 10^9$ with only 7/124 examples demonstrating ill behaviour for such λ . This makes the choice of very large values of λ not attractive at all. Mainly, the analysis shows that choosing $\lambda > 10^9$ could cause the algorithm to diverge. Hence, based on our method, selecting $\lambda \leq 10^7$ seems to be a very safe choice for our BOLIB test set. This is useful for the choice of fixed λ as we can choose values smaller than 10^7 . For the case of varying λ the values are controlled by the stopping criteria proposed in the next section. The complete results on the values of λ_{ill} for each example will be presented in Table 2 below.

Interestingly, 34 out of the 124 test problems do not demonstrate any ill behaviour signs even if we run the algorithm for 1,000 iterations with $\lambda = 0.5 * 1.05^{\text{iter}}$. A potential reason for this could just be that the parameter λ does not get sufficiently large after 1,000 iterations to cause problems for these problems. It could also be that the eigenvalues of the Hessian are not affected by large values of λ for these examples. Also, there is a possibility that elements of the Hessian involved in (2.19) do not depend on λ at all, as for 20/34 problems where the function g is linear in (x, y) or not

present in these problems. Next, we focus on assessing what magnitudes of the penalty parameter λ seem to lead to the best performance for our method.

3.2 Do the values of λ really need to be large?

It is clear from the previous subsection that to reduce the chances for Algorithm 2.6 to diverge or exhibit some ill-behaviour, we should approximately select $\lambda < 10^7$. However, it is still unclear whether only large values of λ would be sensible to ensure a good behaviour of the algorithm. In other words, it is important to know whether relatively small values of λ can lead to a good behaviour for Algorithm 2.6. To assess this, we attempt here to identify inflection points, i.e., values of k where we have $\|\Upsilon_{\mu}^{\lambda_k}(z^k)\| < \epsilon$ as λ_k varies increasingly as described in the introduction of this section. We would then record the value of λ_k at these points.

Ideally, we want to get the threshold $\bar{\lambda}$ such that solution is retained for all $\lambda > \bar{\lambda}$ in the sense of Theorem 2.2. To proceed, we extract the information on the final $\text{Error}^* := \|\Upsilon^{\lambda}(z^*)\|$ for each example from [26] and then rerun the algorithm with varying penalty parameter $\lambda_k := 0.5 * 1.05^k$ with new stopping criterion (i.e., $\text{Error}_k \leq 1.1\text{Error}^*$) while relaxing all of the stopping criteria (see details in Section 4). This way, we stop once we observe $\text{Error}_k := \|\Upsilon^{\lambda_k}(z)\|$ close to the Error^* that we obtained in our experiments [26]. It is worth noting that it would make sense to test only 72/117 (61.54%) of examples, for which algorithm performed well and produced a good solution. This approach can be thought of as finding the inflection point. For instance, if we have an algorithm running as below, we want to stop at the inflection point after 125-130 iterations. The illustration of how we aim to stop at the inflection point is presented in Figure 2 (a)–(b) below, where we have $\|\Upsilon^{\lambda}(z)\|$ on the y-axis and iterations on the x-axis.

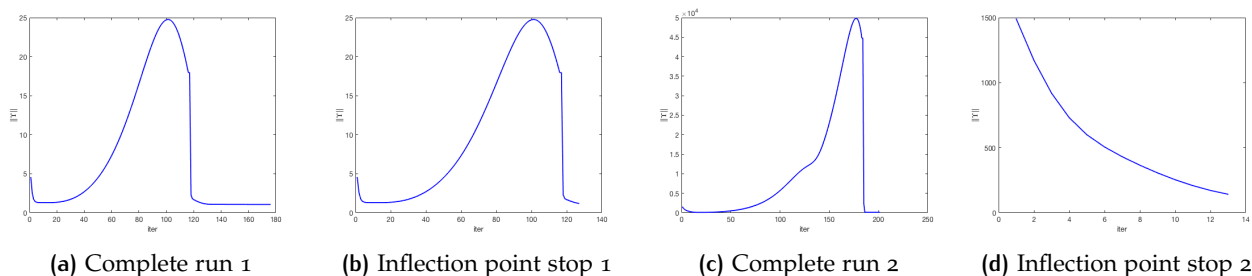


Figure 2: (a) and (b) illustrate the inflection point identification for Example AllendeStill2013 and similarly (c) and (d) correspond to Example Aneta12009; the examples are taken from BOLIB [33].

For some of the examples, we obtained a better Error than initial Error^* . For these cases, we stopped very early as $\text{Error}_k \leq 1.1\text{Error}^*$ was typically met at an early iteration k (where λ was still small), as demonstrated in Figure 2(c)–(d) above. This demonstrates the disadvantage of λ being an increasing sequence. If the algorithm makes many iterations, the parameter λ keeps increasing without possibility to go back to the smaller values. It turns out that for some examples the smaller value of λ was as good as large values, or even better to recover a solution. This further justifies the choice to start from the small value of λ , that is $\lambda_0 < 1$ and increase it slowly.

With the setting to stop whenever $\text{Error}_k \leq 1.1\text{Error}^*$, we often stop very early. Hence, we do not always get $\bar{\lambda}$ that represents the inflection point, which we aimed to get; cf. Figure 2(c)–(d), where (c) has a scale of 10^4 on the y-axis. Although, we can clearly see from Figure 2 (c) that inflection point lies around 190-200 iterations, where the value of λ is much bigger. It is clear that we stopped earlier due to having small value of Error after 12-45 iterations. It was observed that such scenario is typical for the examples in the considered test set. For this reason, we want to introduce λ^* as the *large threshold* of λ . This value will be used to represent the value of the penalty parameter at the inflection point, where solution starts to be recovered for large values of λ ($\lambda > 6.02$), while we also record $\bar{\lambda}$ as the first (smallest) value of λ for which good solution was obtained. For instance, with λ defined as $\lambda := 0.5 \times 1.05^k$ in Figure 2 (c) we have $\bar{\lambda} = 0.5 \times 1.05^{12}$ and $\lambda^* = 0.5 \times 1.05^{190}$ as we

obtain good solution for small λ after 12 iterations and for large λ after 190 iterations. We shall note that the value $\lambda > 6.02$ is the value of the penalty parameter after we make at least 50 iterations as for the case with varying λ we have $\lambda = 0.5 \times 1.05^{51} = 6.02$.

The complete results of detecting $\bar{\lambda}$ and λ^* is presented in Table 2 below. It was observed that the behaviour of the method follows the same pattern for the majority of the examples. Typically, we get a good solution retaining for a few small values of λ , then value of the Error blows up and takes some iterations to start decreasing, coming back to obtaining and retaining a good solution for large values of λ . Such pattern is clearly demonstrated in Figure 2 (c). Such a behaviour is interesting as usually, only large values of penalty parameters are used in practice [4, 23], as intuitively suggested by Theorem 2.2. As mentioned in [13], some methods require penalty parameters to increase to infinity to obtain convergence.

Table 2: Ill behaviour and two thresholds for λ

Problem number	Problem name	iter_{ill}	λ_{ill}	$\bar{\lambda}$	λ^*	iter for λ^*
1	AiyoshiShimizu1984Ex2	495	1.54e+10	9.92×10^4	9.92×10^4	250
2	AllendeStill2013	Not observed	NA	245	245	127
3	AnEtal2009	Not observed	NA	-	-	-
4	Bard1988Ex1	520	5.22e+10	0.608	245	127
5	Bard1988Ex2	536	1.14e+11	-	-	-
6	Bard1988Ex3	Not observed	NA	0.525	54.1	96
7	Bard1991Ex1	Not observed	NA	0.608	183	121
8	BardBook1998	502	2.17e+10	-	-	-
9	CalamaiVicente1994a	476	6.1e+09	0.608	27.3	82
10	CalamaiVicente1994b	490	1.21e+10	0.739	79.9	104
11	CalamaiVicente1994c	490	1.21e+10	-	-	-
12	CalveteGale1999P1	538	1.26e+11	0.525	1.57×10^3	165
13	ClarkWesterberg1990a	520	5.22e+10	-	-	-
14	Colson2002BIPA1	510	3.2e+10	792	792	151
15	Colson2002BIPA2	900	5.88e+18	1.33	1.65×10^3	166
16	Colson2002BIPA3	110	107	-	-	-
17	Colson2002BIPA4	Not observed	NA	0.551	2.68×10^3	176
18	Colson2002BIPA5	550	2.25e+11	-	-	-
19	Dempe1992a	Not observed	NA	-	-	-
20	Dempe1992b	Not observed	NA	0.551	1.29×10^3	161
21	DempeDutta2012Ex24	Not observed	NA	-	-	-
22	DempeDutta2012Ex31	Not observed	NA	0.525	137	115
23	DempeEtal2012	470	4.55e+09	-	-	-
24	DempeFranke2011Ex41	492	1.33e+10	0.704	44.5	92
25	DempeFranke2011Ex42	495	1.54e+10	0.67	54.1	96
26	DempeFranke2014Ex38	495	1.54e+10	0.551	79.9	104
27	DempeLohse2011Ex31a	502	2.17e+10	0.525	6.02	51
28	DempeLohse2011Ex31b	510	3.2e+10	-	-	-
29	DeSilva1978	495	1.54e+10	0.551	69	101
30	FalkLiu1995	440	1.05e+09	1.1×10^4	1.1×10^4	205
31	FloudasEtal2013	505	2.51e+10	0.525	183	121
32	FloudasZlobec1998	510	3.2e+10	-	-	-
33	GumusFloudas2001Ex1	530	8.5e+10	-	-	-
34	GumusFloudas2001Ex3	510	3.2e+10	-	-	-
35	GumusFloudas2001Ex4	502	2.17e+10	-	-	-
36	GumusFloudas2001Ex5	495	1.54e+10	3.04	38.4	89
37	HatzEtal2013	Not observed	NA	0.608	6.02	51
38	HendersonQuandt1958	Not observed	NA	5.64×10^7	5.64×10^7	380
39	HenrionSurowiec2011	Not observed	NA	0.551	6.02	51
40	IshizukaAiyoshi1992a	495	1.54e+10	-	-	-
41	KleniatiAdjiman2014Ex3	445	1.34e+09	-	-	-
42	KleniatiAdjiman2014Ex4	485	9.46e+09	0.739	42.4	91
43	LamparSagrat2017Ex23	Not observed	NA	0.525	137	115
44	LamparSagrat2017Ex31	Not observed	NA	0.525	6.02	51
45	LamparSagrat2017Ex32	Not observed	NA	0.551	284	130
46	LamparSagrat2017Ex33	495	1.54e+10	0.551	102	109
47	LamparSagrat2017Ex35	Not observed	NA	0.525	298	131
48	LucchettiEtal1987	495	1.54e+10	0.525	6.02	51

49	LuDebSinha2016a	Not observed	NA	0.943	6.02	51
50	LuDebSinha2016b	Not observed	NA	0.67	6.02	51
51	LuDebSinha2016c	Not observed	NA	-	-	-
52	LuDebSinha2016d	890	3.61e+18	-	-	-
53	LuDebSinha2016e	900	5.88e+18	-	-	-
54	LuDebSinha2016f	Not observed	NA	-	-	-
55	MacaLHurter1997	Not observed	NA	0.551	6.02	51
56	MirrLees1999	Not observed	NA	0.579	6.02	51
57	MitsosBarton2006Ex38	398	1.36e+08	6.64	7.32	55
58	MitsosBarton2006Ex39	400	1.5e+08	-	-	-
59	MitsosBarton2006Ex310	470	4.55e+09	56.8	56.8	97
60	MitsosBarton2006Ex311	452	1.89e+09	-	-	-
61	MitsosBarton2006Ex312	470	4.55e+09	0.99	6.02	51
62	MitsosBarton2006Ex313	505	2.51e+10	0.579	54.1	96
63	MitsosBarton2006Ex314	460	2.79e+09	0.855	11.9	65
64	MitsosBarton2006Ex315	470	4.55e+09	26	56.8	97
65	MitsosBarton2006Ex316	Not observed	NA	1.33	6.02	51
66	MitsosBarton2006Ex317	420	3.97e+08	3.7	6.02	51
67	MitsosBarton2006Ex318	Not observed	NA	1.04	6.02	51
68	MitsosBarton2006Ex319	398	1.36e+08	-	-	-
69	MitsosBarton2006Ex320	485	9.46e+09	0.943	6.32	52
70	MitsosBarton2006Ex321	452	1.89e+09	1.09	10.3	62
71	MitsosBarton2006Ex322	470	4.55e+09	0.99	20.4	76
72	MitsosBarton2006Ex323	505	2.51e+10	-	-	-
73	MitsosBarton2006Ex324	495	1.54e+10	0.855	6.02	51
74	MitsosBarton2006Ex325	505	2.51e+10	-	-	-
75	MitsosBarton2006Ex326	505	2.51e+10	-	-	-
76	MitsosBarton2006Ex327	475	5.81e+09	1.26	7.32	55
77	MitsosBarton2006Ex328	510	3.2e+10	-	-	-
78	MorganPatrone2006a	500	1.97e+10	0.525	6.02	51
79	MorganPatrone2006b	505	2.51e+10	0.551	6.02	51
80	MorganPatrone2006c	470	4.55e+09	56.8	56.8	97
81	MuuQuy2003Ex1	Not observed	NA	-	-	-
82	MuuQuy2003Ex2	Not observed	NA	-	-	-
83	NieEtal2017Ex34	520	5.22e+10	0.551	118	112
84	NieEtal2017Ex52	Not observed	NA	-	-	-
85	NieEtal2017Ex54	495	1.54e+10	-	-	-
86	NieEtal2017Ex57	850	5.13e+17	-	-	-
87	NieEtal2017Ex58	780	1.69e+16	-	-	-
88	NieEtal2017Ex61	Not observed	NA	-	-	-
89	Outrata1990Ex1a	505	2.51e+10	0.608	192	122
90	Outrata1990Ex1b	510	3.2e+10	0.551	223	125
91	Outrata1990Ex1c	495	1.54e+10	0.99	718	149
92	Outrata1990Ex1d	495	1.54e+10	-	-	-
93	Outrata1990Ex1e	495	1.54e+10	0.855	873	153
94	Outrata1990Ex2a	520	5.22e+10	0.551	245	127
95	Outrata1990Ex2b	470	4.55e+09	1.78	6.02	51
96	Outrata1990Ex2c	510	3.2e+10	0.551	72.5	102
97	Outrata1990Ex2d	520	5.22e+10	-	-	-
98	Outrata1990Ex2e	470	4.55e+09	0.898	6.02	51
99	Outrata1993Ex31	520	5.22e+10	0.551	144	116
100	Outrata1993Ex32	680	1.28e+14	-	-	-
101	Outrata1994Ex31	910	9.58e+18	-	-	-
102	OutrataCervinka2009	505	2.51e+10	-	-	-
103	PaulaviciusEtal2017a	400	1.5e+08	107	107	110
104	PaulaviciusEtal2017b	490	1.21e+10	-	-	-
105	SahinCircic1998Ex2	510	3.2e+10	-	-	-
106	ShimizuAiyoshi1981Ex1	495	1.54e+10	46.7	46.7	93
107	ShimizuAiyoshi1981Ex2	520	5.22e+10	-	-	-
108	ShimizuEtal1997a	910	9.58e+18	0.638	212	124
109	ShimizuEtal1997b	Not observed	NA	-	-	-
110	SinhaMaloDeb2014TP3	Not observed	NA	0.525	651	147
111	SinhaMaloDeb2014TP6	530	8.5e+10	-	-	-
112	SinhaMaloDeb2014TP7	Not observed	NA	$5.75 * 10^{10}$	$5.75 * 10^{10}$	522

113	SinhaMaloDeb2014TP8	520	5.22e+10	-	-	-
114	SinhaMaloDeb2014TP9	505	2.51e+10	-	-	-
115	SinhaMaloDeb2014TP10	480	7.41e+09	-	-	-
116	TuyEtal2007	495	1.54e+10	2.9	16.8	72
117	Vogel2002	Not observed	NA	-	-	-
118	WanWangLv2011	520	5.22e+10	-	-	-
119	YeZhu2010Ex42	Not observed	NA	0.814	6.02	51
120	YeZhu2010Ex43	Not observed	NA	4.49	6.02	51
121	Yezza1996Ex31	460	2.79e+09	223	223	125
122	Yezza1996Ex41	490	1.21e+10	0.608	46.7	93
123	Zlobec2001a	360	2.12e+07	-	-	-
124	Zlobec2001b	495	1.54e+10	-	-	-

We are now going to proceed with finding the threshold of λ for which we start getting a solution and retain the value for larger values of λ . We are going to proceed with the technique of finding small threshold $\bar{\lambda}$ and large threshold λ^* , which was briefly discussed earlier. To find the threshold we first extract the value of the final Error^* for each example from [26]. To find $\bar{\lambda}$ we run the algorithm with λ being defined as $\lambda := 0.5 \times 1.05^k$, with the new stopping criteria:

Stop if $\text{Error} \leq 1.1\text{Error}^*$.

Of course, we also relax all of the aforementioned stopping criteria, as Error is the main measure here and we know that desired value of Error exists. We then define $\bar{\lambda} := 0.5 \times 1.05^{\bar{k}}$, where \bar{k} is the number of iterations after stopping whenever we get $\text{Error} < 1.1\text{Error}^*$. For most of the cases we detect $\bar{\lambda}$ early due to a good solution after the first few iterations in the same manner as shown in Figure 2 (c). Since for many examples we satisfy condition $\text{Error} \leq 1.1\text{Error}^*$ for early iterations (before the inflection point is achieved), we further introduce λ^* , *the large threshold* λ . The value of λ^* will be used to represent the value of the penalty parameter at the inflection point, where solution starts to be recovered for large values of λ . This will be obtained in the same way as $\bar{\lambda}$ with the only difference that we additionally impose the condition to stop after at least 50 iterations. To obtain λ^* we run the algorithm with $\lambda := 0.5 \times 1.05^k$ and the following stopping criteria:

Stop if $\text{Error} \leq 1.1\text{Error}^*$ & $\text{iter} > 50$.

Then the large threshold is defined as $\lambda^* := 0.5 \times 1.05^{k^*}$, where k^* is the number of iterations after stopping whenever we get $\text{Error} < 1.1\text{Error}^*$ and $k > 50$. This way $\bar{\lambda}$ represents the first (smallest) value of λ for which good solution was obtained, while λ^* represent the actual threshold after which solution is retained for large values of λ . The demonstration of stopping at the inflection point for large threshold λ^* was shown in Figure 2 (b) and for small threshold $\bar{\lambda}$ in Figure 2 (d).

It makes sense to test only the examples where algorithm performed well and produced a good solution. For the rest of the examples, the value of Error^* would not make sense as the measure to stop, and we do not obtain good solutions for these examples by the algorithm anyway. From the optimistic perspective we could consider recovered solutions to be the solutions for which the optimal value of upper-level objective was recovered with some prescribed tolerance. Taking the tolerance of 20%, the total amount of recovered solutions by the method with varying λ is 72/117 (61.54%) for the cases where solution was reported in BOLIB [33]. This result will be shown in more details in Section 4.2. Let us look at the thresholds $\bar{\lambda}$ and λ^* for these examples in the figure below, where the value of λ is shown on y – axis and the example numbers on the x – axis.

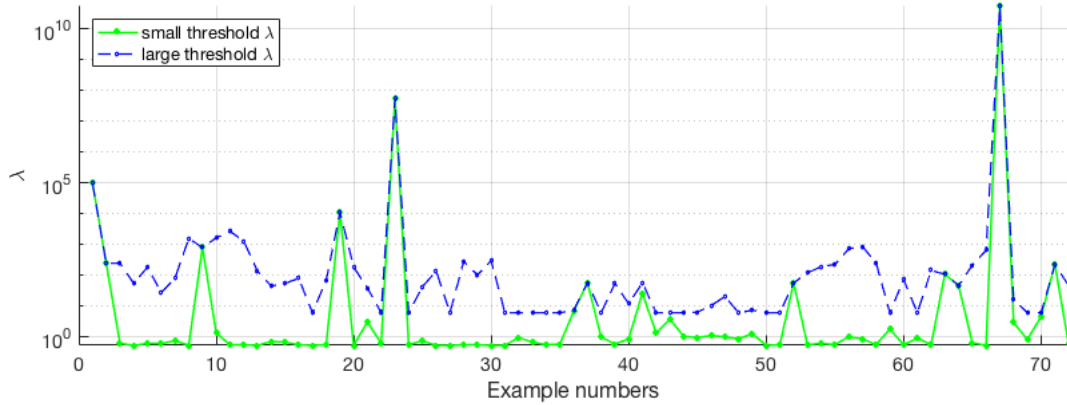


Figure 3: Small threshold $\bar{\lambda}$ and large threshold λ^* for examples with good solutions

For 68/72 problems, we observe that large threshold of the penalty parameter has the value $\lambda^* \leq 10^4$, which shows that we usually can obtain a good solution before 205 iterations. This further justifies stopping algorithm after 200 iterations (see next section) if there is no significant step to step improvement. As for the main outcome of Figure 3, we observe that small threshold is smaller than large threshold ($\bar{\lambda} < \lambda^*$) for 59/72 (82%) problems. This clearly shows that for majority of the problems, for which we recover solution with $\lambda := 0.5 \times 1.05^k$, we obtain a good solution for small λ as well as for large λ . This demonstrates that small λ could in principle be good for the method. For the rest 18% of the problems we have $\bar{\lambda} = \lambda^*$, meaning that good solution was not obtained for $\lambda < 6$ for these examples. This also means that we typically obtain good solution for small values of λ and for large values of λ , but not for the medium values ($\bar{\lambda} < \lambda < \lambda^*$).

As for the general observations of Figure 3, for 42/72 (58.33%) examples we observed that the large threshold λ^* is located somewhere in between 90 – 176 iterations with $40 < \lambda^* < 2680$. For 7/72 problems threshold is in the range $6.02 < \lambda^* \leq 40$, and for only 4/72 problems $\lambda^* > 1.1 \times 10^4$. Once again, this justifies that typically λ does not need to be large. It also suggests the optimal values of λ for the tested examples, at least for our solution method. We observe that for 19/72 problems we get $\lambda^* = 51$, which should be treated carefully as this could mean that the inflection point could possibly be achieved before 50 for these examples. Nevertheless, $\lambda^* = 6.02$ is still a good value of penalty parameter for these examples as solutions are retained for $\lambda > \lambda^*$.

As going to be observed in Section 4 we could actually argue that smaller values of λ work better for our method not only for varying λ but also for fixed λ . Together with the fact that we often have the behaviour as demonstrated in Figure 2 (c), it follows that small λ could be more attractive for the method we implement. We even get better values of Error and better solutions for small values of λ for some examples. Hence we draw the conclusion that small values of λ can generate good solutions. Since it is typical to use large values of λ for other penalization methods (e.g. in [3, 4, 23]), it is interesting what could be the reasons that small λ worked better for our case. This could be due to the specific nature of the method, or due to the fact that we do not do full penalization in the usual sense. Other reason could come from the structure of the problems in the test set. The exact reason of why such behaviour was observed remains an open question. What is important is that this could possibly be the case that small values of λ would be good for some other penalty methods and optimization problems of different nature. This result contradicts typical choice of large penalty parameter for general penalization methods for optimization problems. As the conclusion for our framework, we can claim that for our method λ needs not to be large.

3.3 Partially calm examples

Intuitively, one would think that for partially calm examples, Algorithm 2.6 would behave well, in the sense that varying λ increasingly would lead to a good convergence behaviour. To show that it is not necessarily the case, we start by considering the following result identifying a class of bilevel program of the form (1.1) that is automatically partially calm.

Theorem 3.1 ([19]). Consider a bilevel program (1.1), where G is independent of y and the lower-level optimal solution map is defined as follows, with $c \in \mathbb{R}^m$, $d \in \mathbb{R}^p$, $B \in \mathbb{R}^{p \times m}$, and $A : \mathbb{R}^n \rightarrow \mathbb{R}^p$:

$$S(x) := \arg \min_y \{c^T y \mid A(x) + By \leq d\}. \quad (3.1)$$

In this case, problem (1.1) is partially calm at any of its local optimal solutions.

Examples 8, 40, 43, 45, 46, 188, and 123 in the BOLIB library (see Table 2) are of the form described in this result. The expectation is that these examples will follow the pattern of retaining solution after some threshold, that is for $\lambda > \lambda^*$, as they fit the theoretical structure behind the penalty approach as described in Theorem 2.2. Note that all of these examples follow the pattern shown in Figure 1(a). However, if we relax the stopping criteria used to mitigate the effects of ill-conditioning, as discussed in the previous two subsections, varying λ for 1000 iterations for these seven partially calm examples leads to the 3 typical scenarios demonstrated in Figure 4.

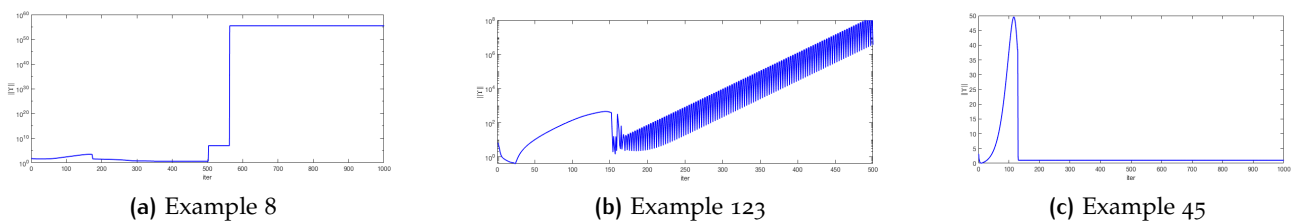


Figure 4: (a) and (b) are obtained for 1000 iterations, while (c) is based on 500 iterations.

In the first case of Figure 4, we can clearly see the algorithm is performing well, retaining the solution for the number of iteration, but then blows up at one point (after 500 iterations) and never goes back to reasonable solution values. Examples 40 and 46 also follow this pattern. Example 123 (second picture in Figure 4) shows a slightly different picture, where the zig-zagging pattern is observed. Algorithm 2.6 blows up at some point and starts zig-zagging away from the solution after obtaining it for a smaller value of λ . Zig-zagging is very common issue in penalty methods and often caused by ill-conditioning [21]. Note that Example 118 exhibits a similar behaviour. This is somewhat similar to scenario 1. However, we put this separately as zig-zagging issue is often referred to as the danger that could be caused by ill-conditioning of a penalty function. The last picture of Figure 4 shows a case where Algorithm 2.6 runs very well without any ill-behaviour observed for all the 1000 iterations. It could be possible that the algorithm could blow up after more iterations if we keep increasing λ . It could also be possible that ill-conditioning does not occur for this example at all, as the Hessian of Υ^λ (2.10) is not affected by λ . Out of the seven BOLIB problems considered here, only examples 43 and 45 follows this pattern.

4 PERFORMANCE COMPARISON FOR THE LEVENBERG-MARQUARDT METHOD UNDER FIXED AND VARYING PARTIAL EXACT PENALTY PARAMETER

Following the discussion from the previous section, two approaches for selecting the penalty parameter λ are tested and compared on the Levenberg-Marquardt method. Recall that for the varying λ case, we define the penalty parameter sequence as $\lambda_k := 0.5 \times 1.05^k$, where k is the iteration number. When fixed values of penalty parameter are considered, ten different values of λ are used for the experiments; i.e., $\lambda \in \{10^6, 10^5, \dots, 10^{-3}\}$. For the fixed values of λ one could choose best λ for each example to see if at least one of the selected values worked well to recover the solution. As in the previous section, the examples used for the experiments are from the Bilevel Optimization LIBrary of Test Problems (BOLIB) [33], which contains 124 nonlinear examples. The experiments are run in MATLAB, version R2016b, on MACI64. Here, we present a summary of the results obtained; more details for each example are reported in the supplementary material [26]. It is important to mention

that algorithm always converges, unlike its Gauss-Newton counterpart studied in [14], where the corresponding algorithm diverges for some values of λ .

4.1 Practical implementation details

For Step 0 of Algorithm 2.6 we set the tolerance to $\epsilon := 10^{-5}$ and the maximum number of iterations to be $K := 1000$. We also choose $\alpha_0 := \|\Upsilon^\lambda(z^0)\|$ (if λ is varying, we set $\lambda := \lambda_0$ here), $\gamma_0 := 1$, $\rho = 0.5$, and $\sigma = 10^{-2}$. The selection of σ is based on the overall performance of the algorithm while the other parameters are standard in the literature.

Starting point. The experiments have shown that the algorithm performs much better if the starting point (x^0, y^0) is feasible. As a default setup, we start with $x^0 = 1_n$ and $y^0 = 1_m$. If the default starting point does not satisfy at least one constraint and algorithm diverges, we choose a feasible starting point; see [26] for such choices. To be more precise, if for some i , $G_i(x^0, y^0) > 0$ or for some j we have $g_j(x^0, y^0) > 0$ and the algorithm does not converge to a reasonable solution, we generate a starting point such that $G_i(x^0, y^0) = 0$ or $g_j(x^0, y^0) = 0$. Subsequently, the Lagrange multipliers are initialised at $u^0 = \max\{0.01, -g(x^0, y^0)\}$, $v^0 = \max\{0.01, -G(x, y)\}$, and $w^0 = 0$.

Smoothing parameter μ . The smoothing process, i.e., the use of the parameter μ , is only applied when necessary (Step 2 and Step 3), where the derivative evaluation for Υ_μ^λ is required. We tried both the case where μ is fixed to be a small constant for all iterations, and the situation where the smoothing parameter is sequence $\mu_k \downarrow 0$. Setting the decreasing sequence to $\mu_k := 0.001/(1.5^k)$, and testing its behaviour and comparing it with a fixed small value ($\mu := 10^{-11}$), in the context of Algorithm 2.6, we observed that both options lead to almost the same results. Hence, we stick to what is theoretically more suitable, that is the smoothing decreasing sequence ; i.e., $\mu_k := 0.001/(1.5^k)$.

Descent direction check and update. For the sake of efficiency, we calculate the direction d^k by solving (2.19) with the Gaussian elimination. Considering the line search in Step 3, if we have $\|\Upsilon^\lambda(z^k + \gamma_k d^k)\|^2 < \|\Upsilon^\lambda(z^k)\|^2 + \sigma \gamma_k \nabla \Upsilon^\lambda(z^k)^\top \Upsilon^\lambda(z^k) d^k$, we redefine $\gamma_k = \gamma_k/2$ and check again. Recall that the Levenberg-Marquardt direction can be interpreted as a combination of Gauss-Newton and steepest descent directions. In fact, if $\alpha_k = 0$ this direction is a Gauss-Newton direction one and when as $\alpha_k \rightarrow \infty$ the direction d^k from (2.19) tends to a steepest descent direction. Hence, if the Levenberg-Marquardt direction is not a descending at some iteration, we give more weight to the steepest descent direction. Hence, when $\|\Upsilon^\lambda(z^k)\| > \|\Upsilon^\lambda(z^{k-1})\|$ setting $\alpha_{k+1} := 10000 \|\Upsilon^\lambda(z^k)\|$ has led to an overall good performance of Algorithm 2.6 for test set used in this paper.

Stopping Criteria. The primary stopping criterion for Algorithm 2.6 is $\|\Upsilon_\mu^\lambda(z^k)\| < \epsilon$, as requested in Step 1. However, robust safeguards are needed to deal with ill-behaviours typically due to the size of the penalty parameter λ . Hence, for the practical implementation of the method, we set $\epsilon = 10^{-5}$ and **stop** if one of the following six conditions is satisfied:

1. $\|\Upsilon^\lambda(z^k)\| < \epsilon$,
2. $|\|\Upsilon^\lambda(z^{k-1})\| - \|\Upsilon^\lambda(z^k)\|| < 10^{-9}$,
3. $|\|\Upsilon^\lambda(z^{k-1})\| - \|\Upsilon^\lambda(z^k)\|| < 10^{-4}$ and iter > 200 ,
4. $\|\Upsilon^\lambda(z^{k-1})\| - \|\Upsilon^\lambda(z^k)\| < 0$ and $\|\Upsilon^\lambda(z^k)\| < 10$ and iter > 175 ,
5. $\|\Upsilon^\lambda(z^k)\| < 10^{-2}$ and iter > 500 ,
6. $\|\Upsilon^\lambda(z^k)\| > 10^2$ and iter > 200 .

The additional stopping criteria is important to ensure that algorithm is not running for too long. The danger of running algorithm for too long is that ill-conditioning could occur. Further, we

typically observe the pattern that we recover solution earlier than algorithm stops. This appears due to the nature of the overdetermined system. We do not know beforehand the tolerance with which we can solve for $\|\Upsilon^\lambda(z)\|$ as Υ^λ is overdetermined system. Hence, it is hard to select ϵ that would fit all examples and allow to solve examples with efficient tolerance. With the stopping criteria defined above we avoid running unnecessary iterations, retaining the obtained solution. To avoid algorithm running for too long and to prevent λ to become too large, we impose additional stopping criterion 3, 5 or 6 above. These criteria are motivated by the behaviour of the algorithm.

For almost all of the examples we observe that after 100-150 iterations we obtain the value reasonably close to the solution but we cannot know beforehand what would be the tolerance of $\|\Upsilon^\lambda(z^k)\|$ to stop. Choosing small ϵ would not always work due to the overdetermined nature of the system being solved. Choosing ϵ too big would lead to worse solutions and possibly not recover some of the solutions. Further, a quick check has shown that ill-conditioning issue typically takes place after 500 iterations for majority of the problems. For these reasons we stop if the improve of the Error value from step to step becomes too small, $|\|\Upsilon^\lambda(z^{k-1})\| - \|\Upsilon^\lambda(z^k)\|| < 10^{-4}$, after the algorithm has performed 200 iterations. Since ill-conditioning is likely to happen after 500 iterations we stop if by that time we obtain a reasonably small Error, $\|\Upsilon^\lambda(z^k)\| < 10^{-2}$. Finally, if it turns out that system cannot be solved with a good tolerance, such that we would obtain a reasonably small value of the Error, we stop if the Error after 200 iterations is big, $\|\Upsilon^\lambda(z^k)\| > 10^2$. This way additional stopping criteria plays the role of safeguard to prevent ill-conditioning and also does not allow the algorithm to keep running for too long once a good solution is obtained.

4.2 Accuracy of the upper-level objective function

Here, we compare the values of the upper-level objective functions at points computed by the Levenberg-Marquardt algorithm with fixed λ and varying λ . For the comparison purpose, we focus our attention only on 117 BOLIB examples [33], as solutions are not known for the other seven problems. To proceed, let \bar{F}_A be the value of upper-level objective function at the point (\bar{x}, \bar{y}) obtained by the algorithm and \bar{F}_K the value of this function at the known best solution point reported in the literature (see corresponding references in [33]). We consider all fixed $\lambda \in \{10^6, 10^5, \dots, 10^{-3}\}$ and varying λ in one graph and present the results in Figure 5 below, where we have the relative error $(\bar{F}_A - \bar{F}_K)/(1 + |\bar{F}_K|)$ on the y-axis and number of examples on the x-axis, starting from 30th example. We further plot the results for the best fixed value of λ . The graph is plotted in the order of increasing error. From the Figure 5 above we can clearly see that much more solutions were recovered for the small values of fixed λ than for large values. For instance with the allowable accuracy error of $\leq 20\%$ we recover solutions for 78.63% for fixed $\lambda \in \{10^{-2}, 10^{-3}\}$, while for $\lambda \in \{10^6, 10^5, 10^4, 10^3, 10^2\}$ we recover at most 40.17% solutions. Interestingly, the worst performance is observed for fixed $\lambda = 100$. With the varying λ we observe that algorithm performed averagely in comparison between large and small fixed values of λ , recovering 59.83% of the solutions with the accuracy error of $\leq 20\%$. It is worth saying that implementing Algorithm 2.6 with $\lambda := 0.5 \times 1.05^k$ still recovers over half of the solutions, which is not too bad. However, fixing λ to be small recovers way more solutions, which shows that varying λ is not the most efficient option for our case.

It was further observed that for some examples only $\lambda \geq 10^3$ performed well, while for others small values ($\lambda < 1$) showed good performance. If we were able to pick best fixed λ for each example, we would obtain negligible (less than 10%) error for upper-level objective function for 85.47% of the tested problems. With the accuracy error of $\leq 25\%$ our algorithm recovered solutions for 88.9% of the problems for the best fixed λ and for 61.54% with varying λ . This means that if one can choose the best fixed λ from the set of different values, fixing λ is much more attractive for the algorithm. It was further observed that for some examples only $\lambda \geq 10^3$ performed well, while for others small values of λ ($\lambda < 1$) showed good performance. However, if one does not have a way to choose the best value or a set of potential values cannot be constructed efficiently for certain problems, varying λ could be a better option to choose. Nevertheless, for the test set of small problems from BOLIB [33], fixing λ to be small performed much better than varying λ as increasing sequence. Further, if

one could run the algorithm for all fixed λ and was able to choose the best one, the algorithm with fixed λ performs extremely well compared to varying λ . In other words, algorithm almost always finds a good solution for at least one value of fixed $\lambda \in \{10^6, 10^5, \dots, 10^{-3}\}$.

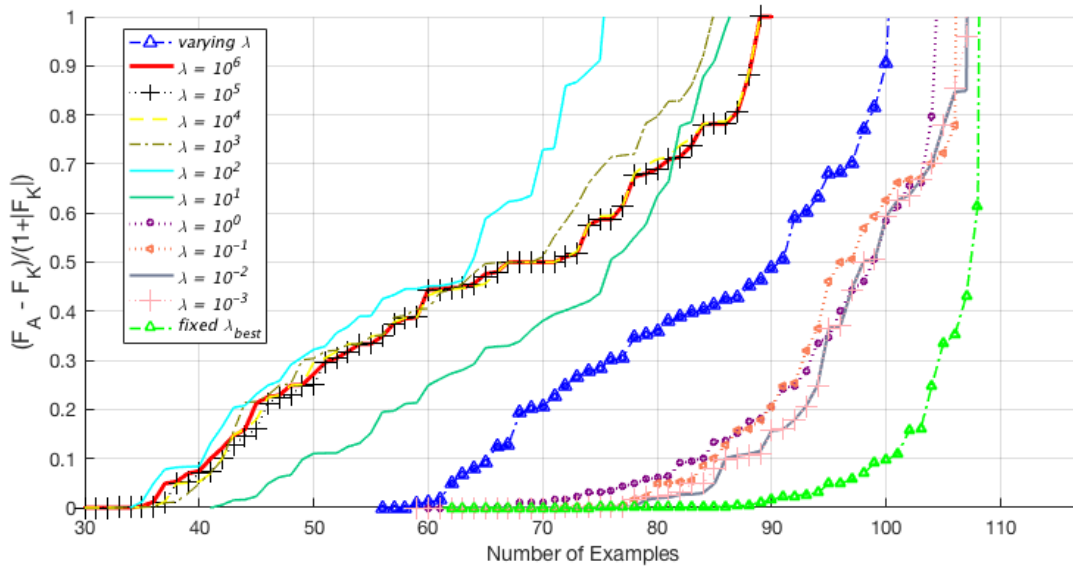


Figure 5: Error of the upper-level objective value for examples with known solutions

4.3 Feasibility check

Considering the structure of the feasible set of problem (1.2), it is critical to check whether the points computed by our algorithms satisfy the value function constraint $f(x, y) \leq \varphi(x)$, as it is not explicitly included in the expression of Υ^λ (2.10). If the lower-level problem is convex in y and a solution generated by our algorithms satisfies (2.5) and (2.8), then it will verify the value function constraint. Conversely, to guaranty that a point (x, y) such that $y \in S(x)$ satisfies (2.5) and (2.8), a constraint qualification (CQ) is necessary. Note that conditions (2.5) and (2.8) are incorporated in the stopping criterion of Algorithm 2.6. To check whether the points obtained are feasible, we first identify the BOLIB examples, where the lower-level problem is convex w.r.t. y . As shown in [14] it turns out that a significant number of test examples have linear lower-level constraints. For these examples, the lower-level convexity is automatically satisfied. We detect 49 examples for which some of these assumptions are not satisfied, that is the problems having non-convex lower-level objective or some of the lower-level constraints being nonconvex. For these examples, we compare the obtained solutions with the known ones from the literature. Let f_A stand for $f(\bar{x}, \bar{y})$ obtained by one of the tested algorithms and f_K to be the known optimal value of lower-level objective function. In the graph below we have the lower-level relative error, $(f_A - f_K)/(1 + |f_K|)$, on the y-axis, where the error is plotted in increasing order. In Figure 6 below we present results for all fixed $\lambda \in \{10^6, 10^5, \dots, 10^{-3}\}$ as well as varying λ defined as $\lambda := 0.5 \times 1.05^k$.

From the Figure 6 above we can see that for 20 problems the relative error of lower-level objective is negligible ($< 5\%$) for all values of fixed λ and varying λ . We have seen that convexity and a CQ hold for the lower-level hold for 74 test examples. We consider solutions for these problems to be feasible for the lower-level problem. Taking satisfying feasibility error to be $< 20\%$ and using information from the graph above, we claim that feasibility is satisfied for at most 100 (80.65%) problems for fixed $\lambda \in \{10^6, 10^5, 10^4, 10^3\}$, for 101 – 104 (81.45 – 83.87%) problems for $\lambda \in \{10^3, 10^2, 10^1, 10^0, 10^{-1}\}$ and for 106 (85.48%) problems for $\lambda \in \{10^{-2}, 10^{-3}\}$. We further observe that feasibility is satisfied for 101 (81.4%) problems for varying λ . Considering we could choose best fixed λ for each of the examples, we could also claim that feasibility is satisfied for 108 (87.1%) problems for best fixed λ . From Figure 6 we note that slightly better feasibility was observed for smaller values of fixed λ than

for the big ones and that varying λ has shown average performance between these magnitudes in terms of the feasibility.

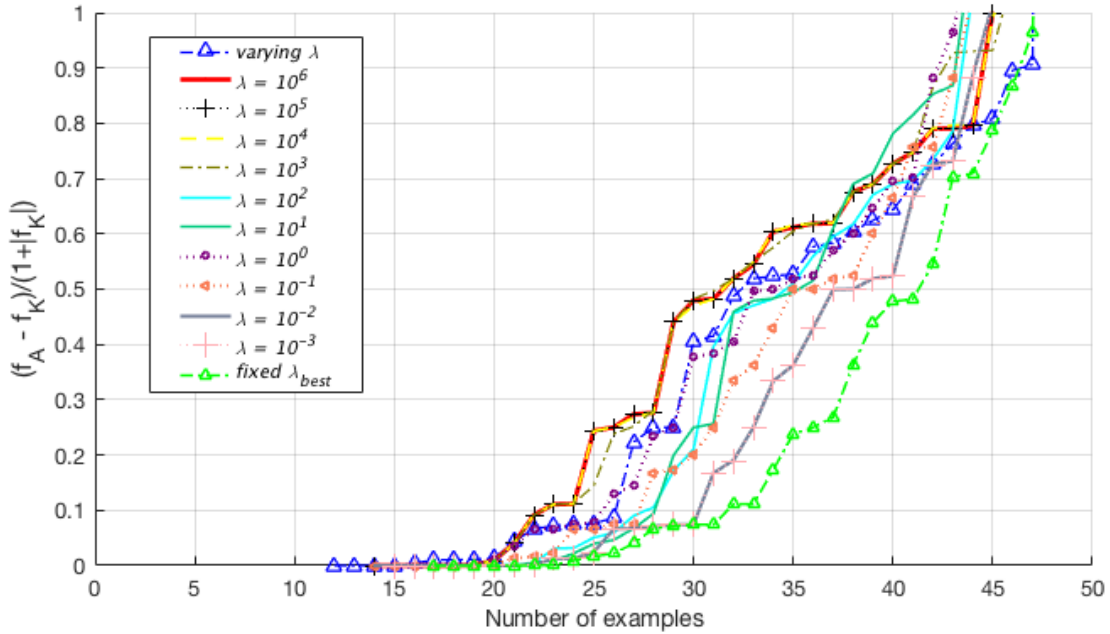


Figure 6: Feasibility error for the lower-level problem in increasing order

4.4 Experimental order of convergence

Recall that the experimental order of convergence (EOC) is defined by

$$\text{EOC} := \max \left\{ \frac{\log \|\gamma^\lambda(z^{K-1})\|}{\log \|\gamma^\lambda(z^{K-2})\|}, \frac{\log \|\gamma^\lambda(z^K)\|}{\log \|\gamma^\lambda(z^{K-1})\|} \right\},$$

where K is the number of the last iteration [12]. If $K = 1$, no EOC will be calculated ($\text{EOC} = \infty$). EOC is important to estimate the local behaviour of the algorithm and to show whether this practical convergence reflects the theoretical convergence result stated earlier. Let us consider EOC for fixed $\lambda \in \{10^6, \dots, 10^{-3}\}$ and for varying λ ($\lambda = 0.5 \times 1.05^k$) in Figure 7 below.

It is clear from this picture that for most of the examples our method has shown linear experimental convergence. This is slightly below the quadratic convergence established by Theorem 2.7. It is however important to note that the method always converges to a number, although sometimes the output might not be the optimal point for the problem. There are a few examples that shown better convergence for each value of λ , with the best ones being $\lambda \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^6\}$ as seen in the figure above. These fixed values have shown slightly better EOC performance than varying λ . Varying λ showed slightly better convergence than fixed $\lambda \in \{10^0, 10^1, 10^2, 10^3, 10^4, 10^5\}$. EOC bigger than 1.2 has been obtained for less than 5 (4.03 %) examples for fixed $\lambda \in \{10^0, 10^1, 10^2, 10^3, 10^4, 10^5\}$, while varying λ showed such EOC for 11 (8.87%) examples. Fixed $\lambda = 10^6$ has shown almost the same result as varying λ with rate of convergence greater than 1.2 for 12 (9.67%) examples, while $\lambda = 10^{-1}$ has demonstrated such EOC for 14 (11.29%) examples and $\lambda \in \{10^{-2}, 10^{-3}\}$ for 17 (13.71 %) examples. Finally, in the graph above, we can see that for all values of λ only a few ($\leq 4/124$) examples have worse than linear convergence.

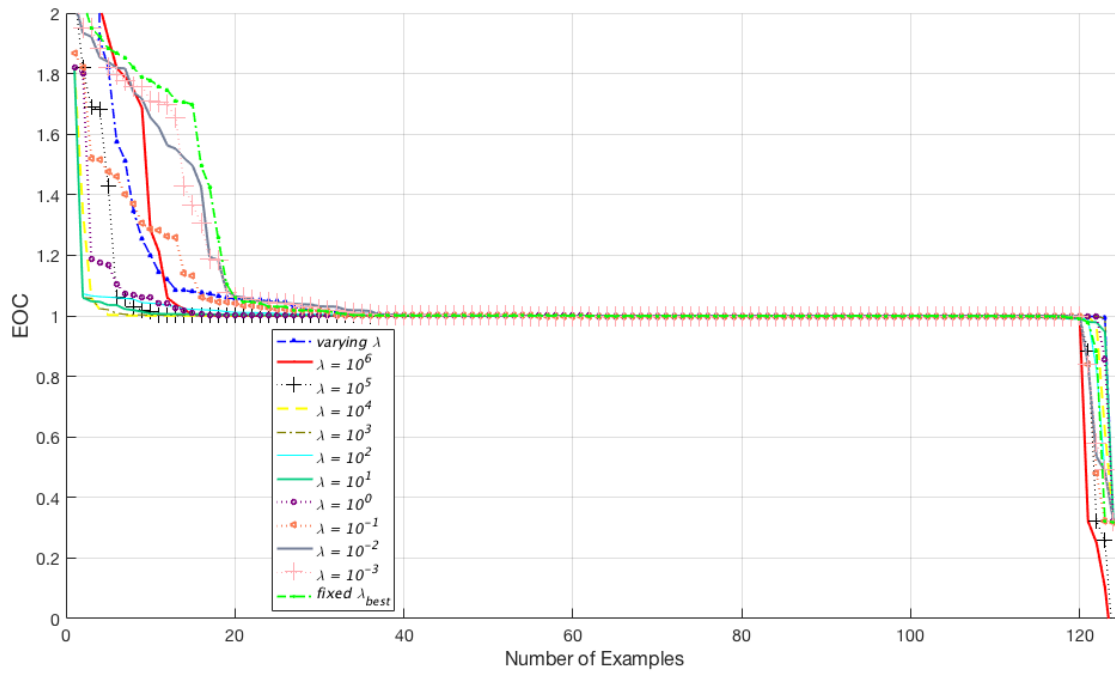


Figure 7: Observed EOC at the last iterations for all examples (in decreasing order)

4.5 Line search stepsize

Let us now look at the line search stepsize, γ_k , at the last step of the algorithm for each example. Consider all fixed λ and varying λ in Figure 8 below. This is quite important to know two things. Firstly, how often line search was used at the last iteration, that is how often implementation of line search was clearly important. Secondly, as main convergence results are for the pure method this would be demonstrative to note how often the pure (full) step was made at the last iteration. This can then be compared with the experimental convergence results in the previous subsection, namely with Figure 7.

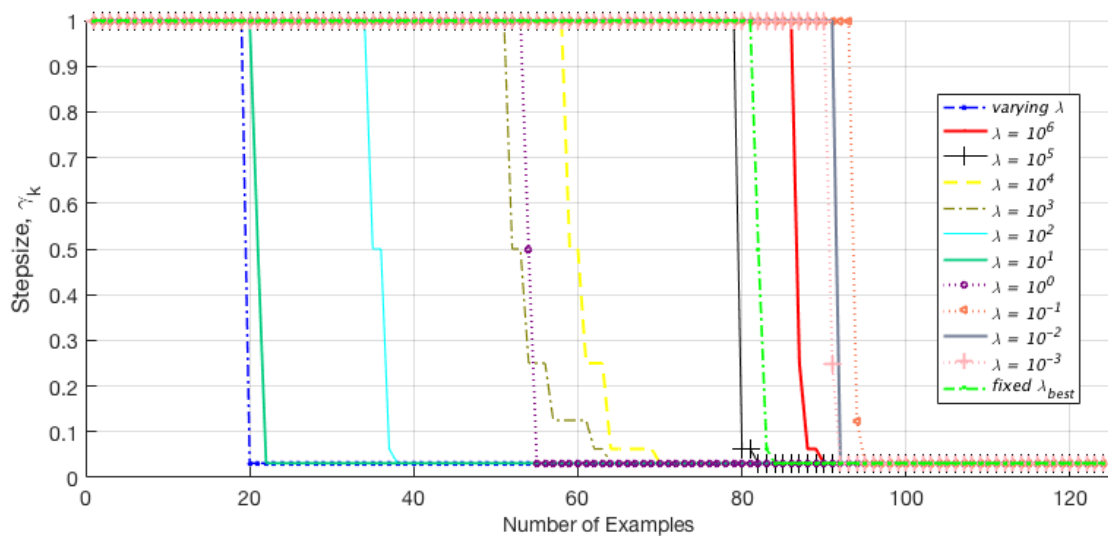


Figure 8: Stepsize made at the last iteration for all examples (in decreasing order)

In the figure above whenever stepsize on the y-axis is equal to 1 it means the full step was made at the last iteration. For these cases the convergence results shown in Theorem 2.7 could be considered valid. From the graphs above we observe that stepsize at the last iteration was rather $\gamma_k = 1$ or $\gamma_k <$

0.05. We observe that for varying λ algorithm would typically do a small step at the last iteration. It seems that algorithm with varying λ benefits more from line search technique than the algorithm with fixed λ . Possibly, pure Levenberg-Marquardt method with varying λ would not converge for most of the problems. Interestingly, for fixed values of λ stepsize was $\gamma_k < 0.05$ at the last iteration much more often for the values of λ that showed worse performance in terms of recovering solutions (i.e. $\lambda \in \{10^1, 10^2\}$). We also observe that for medium values of λ ($\lambda \in \{10^4, 10^3, 10^2, 10^1, 10^0\}$) full stepsize was made for less than half of the examples. For large values $\lambda \in \{10^5, 10^5\}$ full step was made for 63.71% and 70.16% of the problems respectively. Further on, small values of λ for which more solutions were recovered would do the full step at the last iteration for most of the examples. For instance, with $\lambda = 10^{-3}$ and $\lambda = 10^{-2}$ full step was made at the last iteration for 73.39% of the problems, while for $\lambda = 10^{-1}$ full step was made for 75.81% of the problems. In terms of the fixed λ_{best} it is interesting to observe that full step was used only for 82/124 (66.13%) of the problems, meaning that for a third of the problems linesearch was implemented in the last step for the best tested value of λ . This also coincides with the results of Figure 7 where with smaller values of λ algorithm has shown faster than linear convergence for more examples than for big values of λ . This is likely to be the case that small steps were made in the other instances due to the non-efficient direction of the method at the last iteration.

5 FINAL COMMENTS

We introduce a smoothing Levenberg-Marquardt method to solve LLVF-based optimality conditions for optimistic bilevel programs. Since these optimality conditions are parameterized by a certain λ , its selection is carefully studied in light of the BOLIB library test set. Surprisingly, small values of this partial exact penalty parameter showed a good behaviour for our method, as they generally performed very well, although classical literature on exact penalization usually suggest large values. However, relatively medium values of λ did not perform as well. Furthermore, as both the varying and fixed scenarios were considered for λ , it was observed that both approaches showed a linear experimental order of convergence for most of the examples. Average CPU time for all fixed λ is 0.243 seconds, while it is 0.193 seconds for λ_{best} and 0.525 seconds for varying λ . The algorithm with varying λ turns out to be more than twice slower than for fixed λ . However, running it for all fixed values of λ can take way more time than the case where the value of λ is varying.

REFERENCES

- [1] G.B. Allende and G. Still, *Solving bi-level programs with the KKT-approach*, Mathematical Programming 131:37-48 (2012)
- [2] R. Behling, A. Fischer, G. Haeser, A.Ramos and K. Schönefeld, *On the constrained error bound condition and the projected Levenberg-Marquardt method*, Optimization 66(8): 1397-1411 (2016)
- [3] D.P. Bertsekas, *Constrained optimization and Lagrange multiplier methods*, Academic Press (1982)
- [4] J.V. Burke, *An exact penalization viewpoint of constrained optimization*, SIAM journal on control and optimization 29(4): 968-998 (1991)
- [5] S. Dempe and J. Dutta, *Is bilevel programming a special case of mathematical programming with equilibrium constraints?* Mathematical Programming 131:37-48 (2012)
- [6] S. Dempe, J. Dutta, and B.S. Mordukhovich, *New necessary optimality conditions in optimistic bilevel programming*, Optimization 56 (5-6):577-604 (2007)
- [7] S. Dempe and A.B. Zemkoho (eds.), *Bilevel optimization: advances and next challenges*, Springer (2020)
- [8] S. Dempe and A.B. Zemkoho, *The bilevel programming problem: reformulations, constraint qualification and optimality conditions*, Mathematical Programming 138:447-473 (2013)
- [9] S. Dempe and A.B. Zemkoho, *The generalized Mangasarian-Fromowitz constraint qualification and optimality conditions for bilevel programs*, Journal of Optimization Theory and Applications 148(1):46-68 (2011)

- [10] J.Y. Fan and Y.X. Yuan, *On the quadratic convergence of the Levenberg-Marquardt method without nonsingularity assumption*, *Computing* 74:23-39 (2005)
- [11] A. Fischer, *A special Newton-type optimization method*, *Optimization* 24(3):269-284 (1992)
- [12] A. Fischer, A.B. Zemkoho, and S. Zhou, *Semismooth Newton-type method for bilevel optimization: global convergence and extensive numerical experiments*, arXiv:1912.07079 (2019)
- [13] R. Fletcher, *An ideal penalty function for constrained optimization*, *IMA Journal of Applied Mathematics* 15:319-342 (1975)
- [14] J. Fliege, A. Tin, and A.B. Zemkoho, *Gauss-Newton-type methods for bilevel optimization*, *Computational Optimization and Applications*, doi.org/10.1007/s10589-020-00254-3 (2021)
- [15] J. Herskovits, M.T. Filho, and A. Leontiev, *An interior point technique for solving bilevel programming problems*, *Optimization and Engineering* 14(3):381-394 (2013)
- [16] C. Kanzow, *Some noninterior continuation methods for linear complementarity problems*, *SIAM Journal on Matrix Analysis and Applications* 17(4):851-868 (1996)
- [17] L. Lampariello and S. Sagratella, *Numerically tractable optimistic bilevel problems*, *Computational Optimization and Applications* 76(2):277-303 (2020)
- [18] G.-H. Lin, M. Xu, and J.J. Ye, *On solving simple bilevel programs with a nonconvex lower level program*, *Mathematical Programming* 144(1-2):277-305 (2014)
- [19] P. Mehrlitz, L.I. Minchenko, and A.B. Zemkoho, *A note on partial calmness for bilevel optimization problems with linear structures at the lower level*, *Optimization letters*, doi.org/10.1007/s11590-020-01636-6 (2020)
- [20] A. Mitsos, P. Lemonidis, and P.I. Barton, *Global solution of bilevel programs with a nonconvex inner program*, *Journal of Global Optimization* 42(4):475-513 (2008)
- [21] J. Nocedal and S.J. Wright, *Numerical optimization*, Springer (1999)
- [22] R. Paulavicius, J. Gao, P. Kleniati, and C.S. Adjiman, *BASBL: Branch-and-sandwich bilevel solver. Implementation and computational study with the BASBLib test sets*, *Computers & Chemical Engineering* 132 (2020)
- [23] G. Di Pillo and L. Grippo, *Exact penalty functions in constrained optimization*, *SIAM journal on control and optimization* 27 (6): 1333-1360, (1989)
- [24] S. Pineda, H. Bylling, and J.M. Morales, *Efficiently solving linear bilevel programming problems using off-the-shelf optimization software*, *Optimization and Engineering* 19(1):187-211 (2018)
- [25] W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery, *Numerical recipes in C: the art of scientific computing (2nd Ed)*, Cambridge University Press (1992)
- [26] A. Tin and A.B. Zemkoho, *Supplementary material for "Levenberg-Marquardt method for bilevel optimization"*, School of Mathematical Sciences, University of Southampton, UK (2020)
- [27] W. Wiesemann, A. Tsoukalas, P. Kleniati, and B. Rustem, *Pessimistic bilevel optimization*, *SIAM Journal on Optimization* 23(1):353-380 (2013)
- [28] M. Xu and J.J. Ye, *A smoothing augmented lagrangian method for solving simple bilevel programs*, *Computational Optimization and Applications* 59(1-2):353-377 (2014)
- [29] N. Yamashita and M. Fukushima, *On the rate of convergence of the Levenberg-Marquardt method*, *Computing* 15:237-249 (2001)
- [30] J.J. Ye and D.L. Zhu, *Optimality conditions for bilevel programming problems*, *Optimization* 33:9-27 (1995)
- [31] A.B. Zemkoho and S. Zhou, *Theoretical and numerical comparison of the Karush–Kuhn–Tucker and value function reformulations in bilevel optimization*, *Computational Optimization and Applications* doi.org/10.1007/s10589-020-00250-7 (2021)
- [32] A.B. Zemkoho, *Bilevel programming: reformulations, regularity and stationarity*, PhD thesis, Department of Mathematics and Computer Science, TU Bergakademie Freiberg, Freiberg, Germany (2012)
- [33] S. Zhou, A.B. Zemkoho, and A. Tin, *BOLIB: Bilevel Optimization LIBrary of test problems*, In: S. Dempe and A.B. Zemkoho (eds.), *Bilevel optimization: advances and next challenges*, Springer (2020)