
A Trust-Region Interior-Point Method for Bilevel Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Bilevel optimization (BLO) is popular in designing learning tasks in hyperparameter
2 tuning, meta learning, reinforcement learning and adversarial learning. Most
3 existing methods for solving BLO are gradient based methods, which often need to
4 solve a lower level problem approximately to obtain an approximate hyper-gradient.
5 In this paper, we propose the first second-order interior-point method (IPM) based
6 on value function approach for solving BLO, i.e., the Bilevel Trust-Region Interior-
7 Point Method (BTRIPM). As the value function reformulation is nonconvex, we
8 adopt the trust-region method to solve the log-barrier subproblem. Like IPMs in
9 nonlinear optimization, the BTRIPM admits empirical rapid convergence. We
10 theoretically prove convergence and rate of convergence of the proposed method
11 under mild conditions that are widely used in nonlinear IPM or BLO community.
12 Experiments on a toy example and hyperparameter tuning with real-world datasets
13 demonstrate the efficiency and accuracy of the proposed method over existing
14 first-order methods.

15 1 Introduction

16 Bilevel Optimization (BLO) refers to a type of Optimization problems with hierarchical structures.
17 It is widely applied in practical machine learning models [5, 9, 12], such as hyper-parameter opti-
18 mization [19, 15], meta learning [20, 8], reinforcement learning [27], adversarial learning [3, 2, 25].
19 Generally, a BLO takes the following formulation

$$\min_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathbb{R}^n} F(\mathbf{x}, \mathbf{y}), \text{ s.t. } \mathbf{y} \in \operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}), \quad (1)$$

20 where $F : \mathcal{X} \times \mathbb{R}^n \rightarrow \mathbb{R}$ is called the Upper-Level (UL) objective and $f : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ is the Lower-
21 Level (LL) objective. In the literature [19, 7], the solution set of the LL problem $\operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$
22 is often required to be a singleton, denoted as $\{\mathbf{y}^*(\mathbf{x})\}$, in order that (1) can be reformulated as
23 a single-level optimization problem, i.e. $\min_{\mathbf{x} \in \mathcal{X}} \phi(\mathbf{x}) = F(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))$. We call $\partial_{\mathbf{x}} \phi$ and $\partial_{\mathbf{x}\mathbf{x}}^2 \phi$ the
24 hyper-gradient and hyper-Hessian of (1), respectively.

25 1.1 Related Work

26 Generally, BLO is intrinsically NP hard [1], where the difficulty lies in dealing with the special
27 constraint. In the literature, there are various approaches for solving BLO. To find an approximate
28 hyper-gradient, Explicit Gradient-Based Methods (EGBMs) [7, 8] use dynamics on iterative algo-
29 rithms to solve the LL problem. In this framework, Reverse Hyper-Gradient (RHG) and Forward
30 Hyper-Gradient (FHG) methods identify the hyper-gradient by forward and reverse computation
31 iterations, respectively. To ease the computation, Shaban et al. [22] develops a technique that truncates

the back-propagation process to reduce the computation. Besides, Implicit Gradient-Based Methods (IGBMs) [19, 20, 14] are also prevalent for BLO. Using the first-order optimality condition for LL problem and the chain rule, IGBMs solve a linear system to calculate the hyper-gradient. Even approximately solving the linear system by the Conjugate Gradient (CG) method or the Neumann method as widely used in the literature [19, 14], IGBMs demand huge computation. To avoid the expensive Hessian vector products, Liu et al. [13] recently proposes an algorithm termed Bilevel Value-Function-based Interior-point Method (BVFIM). By approximating the constraint with a series of inequalities based on value functions of the LL constraint, the BVFIM uses the log-penalty function to combine the UL and LL objectives. More specifically, in each step, the BVFIM first approximately minimizes the LL problem w.r.t. \mathbf{y} at current $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ to obtain a value function representation of the LL problem, then minimizes, w.r.t. \mathbf{y} , the objective function penalized by the log barrier of the value function inequality to achieve a smooth approximation of $\phi(\mathbf{x})$, and finally applies gradient descent to update \mathbf{x} using this smooth approximation. Numerical experiments in Liu et al. [13] show BVFIM outperforms existing methods.

However, to the best of our knowledge, there are almost no second-order algorithms considered in the BLO literature or in the machine learning society, though some works in BLO [17, 29] analyzed the second-order optimality conditions. One reason may be that it is very difficult to estimate the hyper-Hessian $\partial_{\mathbf{x}\mathbf{x}}^2 \phi(\mathbf{x})$. As is well known, second-order methods enjoy rapid convergence in general nonlinear optimization. This motivates us to design the algorithm in this paper.

1.2 Our Contributions

In this paper, we propose a Bilevel Trust-Region Interior-Point Method (BTRIPM), which is the first value function based second-order method for BLO. We approximate the LL constraint $\mathbf{y} \in \arg\min_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ by an inequality associated with its value function following value function approaches in the literature [13, 28]. We penalize a relaxation of this value function inequality with the log-barrier penalty and all other inequalities in \mathcal{X} as in the interior-point method (IPM) literature [6]. Unlike Liu et al. [13], we minimize a sequence of penalized problems regarding both \mathbf{x} and \mathbf{y} as decision variables. As the log-barrier problems are possibly nonconvex, we distinguish our algorithm from convex IPMs by applying a trust-region method to solve the log-barrier problems.

The Hessian vector product in our BTRIPM can be computed in a cost dominated by solving a linear system in dimension with n equations. The cost can further be reduced and fast convergence can still be guaranteed in practice if we use the upper level Hessian to approximate the true Hessian when the dimension of LL problem is high. As a second-order method, our algorithm converges faster than first-order methods in the literature [7, 8, 22, 19, 20, 14, 19, 14, 13]. Our BTRIPM needs often several tens of outer iterations to obtain a solution with higher precision and the computational time is less than the-state-of-the-art methods in the experiments.

Moreover, we theoretically prove that the proposed algorithm converges to a strict local optimal solution under mild assumptions. Our proof technique is totally different from the existing first-order methods for BLO. Our technique successfully addresses the well known difficulty in analyzing the value function approach that many usual constraint qualification (CQ) such as the nonsmooth Mangasarian Fromovitz constraint qualification fail at each feasible point (see, e.g., [28]). Specifically, we show that the sequence generated by our algorithm converges to a KKT point of a relaxation of the value function reformulation of (1) that the linear independence constraint qualification (LICQ) holds under minor conditions, and show that the KKT point is a strict local minimum of (1) under some additional minor conditions.

In summary, our contributions are as follows:

- We propose the BTRIPM, the first second-order interior-point method for BLO. The BTRIPM first uses value function approach to rewrite the LL problem as an inequality constraint, and then applies trust-region methods to minimize a sequence of log-barrier problems.
- We are the first to prove the local minimizers of a relaxed value function reformulation of BLO converge to the local minimizer of the original problem. Based on this, we prove that our algorithm converges to a strict local minimizer of (1) under mild conditions.
- Our experiments show that the BTRIPM is faster and more accurate than existing methods in the literature on both toy example and real-world datasets.

Notation. In this paper we consider $\mathcal{X} = \{\mathbf{x} : \mathbf{c}(\mathbf{x}) \geq 0\}$ where $\mathbf{c} : \mathbb{R}^m \rightarrow \mathbb{R}^\kappa$, κ is the number of the inequality constraints.¹ Through the paper, we assume the UL and LL functions F and f , and $c_i, i \in [\kappa]$ are twice continuously differentiable functions in its domain. We define derivatives as follows: for a scalar function $h_1(\mathbf{x})$,

$$\partial h_1(\mathbf{x}) \triangleq \partial_{\mathbf{x}} h_1(\mathbf{x}) = \frac{\partial h_1(\mathbf{x})}{\partial \mathbf{x}} = \left(\frac{\partial h_1(\mathbf{x})}{\partial x_1} \dots \frac{\partial h_1(\mathbf{x})}{\partial x_m} \right)^T \in \mathbb{R}^m,$$

and for a vector function $\mathbf{h}(\mathbf{x}) \in \mathbb{R}^n$,

$$\partial_{\mathbf{x}} \mathbf{h}(\mathbf{x}) = \frac{\partial \mathbf{h}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial h_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial h_1(\mathbf{x})}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_n(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial h_n(\mathbf{x})}{\partial x_m} \end{pmatrix} \in \mathbb{R}^{n \times m}.$$

Following this, we further define

$$\partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) \triangleq \partial_{\mathbf{y}} (\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})) \in \mathbb{R}^{m \times n}.$$

For notational convenience, for $c_i(\mathbf{x})$ as a function of \mathbf{x} , we always use the convention $\partial c_i(\mathbf{x}) = (\partial_{\mathbf{x}} c_i(\mathbf{x})^T, \mathbf{0}^T)^T \in \mathbb{R}^{m+n}$, and if we want to specify the gradient of c_i w.r.t. \mathbf{x} , we use $\partial_{\mathbf{x}} c_i(\mathbf{x})$. Given $\mathbf{x} \in \mathcal{X}$, we always use $\mathbf{z}^*(\mathbf{x})$ to denote an element of $\text{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$, which is a vector function of \mathbf{x} and not unique if $\text{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is not a singleton. Several definitions, technical lemmas and all the proofs defer to the appendix.

2 The Bilevel Trust-region Interior-point Method

In this section, we provide a new algorithm that incorporates the second-order information using nonconvex IPM framework.

To begin with, we reformulate the LL constraint into an inequality constraint based on the value function approach following [13, 28]. Specifically, (1) is equivalent to

$$\min_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathbb{R}^n} F(\mathbf{x}, \mathbf{y}), \text{ s.t. } f(\mathbf{x}, \mathbf{y}) \leq f^*(\mathbf{x}), \quad (2)$$

where $f^*(\mathbf{x}) = \min_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$. However, the constraint is not friendly to IPMs as we always have $f^*(\mathbf{x}) \leq f(\mathbf{x}, \mathbf{y})$ and thus there is no interior point. To avoid this, we introduce a parameter $\mu > 0$ to relax the inequality constraint and obtain

$$\min_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathbb{R}^n} F(\mathbf{x}, \mathbf{y}), \text{ s.t. } f(\mathbf{x}, \mathbf{y}) \leq f^*(\mathbf{x}) + \mu. \quad (3)$$

As $\mathcal{X} = \{\mathbf{x} : \mathbf{c}(\mathbf{x}) \geq 0\}$, (3) is then equivalent to

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n} \quad & F(\mathbf{x}, \mathbf{y}) \\ \text{s.t.} \quad & f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}) \geq 0, \quad c_i(\mathbf{x}) \geq 0, i \in [\kappa], \end{aligned} \quad (4)$$

where $f_{\mu}^*(\mathbf{x}) \triangleq f^*(\mathbf{x}) + \mu$. Using the log-barrier penalty as in IPMs, we obtain the following penalized problem

$$\min_{\mathbf{x}, \mathbf{y}} g(\mathbf{x}, \mathbf{y}), \quad (5)$$

where

$$g(\mathbf{x}, \mathbf{y}) \triangleq F(\mathbf{x}, \mathbf{y}) - \tau \ln(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})) - \tau \sum_{i=1}^{\kappa} \ln(c_i(\mathbf{x})).$$

To achieve a fast convergence in solving the unconstrained problem (5), in contrast to Liu et al. [13] we consider a second-order algorithm rather than use gradient descent. Moreover, as (5) is nonconvex, we cannot use damped Newton's method like usual IPMs. Instead, we apply the trust-region method that ensures monotonic decrease of objective value to solve (5) [4].

In the following we explore the formula of first- and second-order derivatives of $g(\mathbf{x}, \mathbf{y})$.

¹We remark our algorithm allows additional affine constraints. For notational simplicity in the convergence analysis in Section 3, we only consider inequality constraints in this paper.

Table 1: The details of the second-order derivative of $\varphi(x, y)$

Notation	Expression
$\partial_{\mathbf{x}\mathbf{x}}^2 g(\mathbf{x}, \mathbf{y})$	$\partial_{\mathbf{x}\mathbf{x}}^2 F(\mathbf{x}, \mathbf{y}) + \tau \frac{\partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}, \mathbf{y}) - \partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})} + \tau \frac{[\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) - \partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x})][\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) - \partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x})]^T}{(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))^2}$ $+ \tau \sum_{i=1}^{\kappa} [\frac{\partial_{\mathbf{x}} c_i(\mathbf{x}) \partial_{\mathbf{x}} c_i(\mathbf{x})^T}{c_i(\mathbf{x})^2} - \frac{\partial_{\mathbf{x}\mathbf{x}}^2 c_i(\mathbf{x})}{c_i(\mathbf{x})}]$
$\partial_{\mathbf{y}\mathbf{y}}^2 g(\mathbf{x}, \mathbf{y})$	$\partial_{\mathbf{y}\mathbf{y}}^2 F(\mathbf{x}, \mathbf{y}) + \tau \frac{\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{y})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})} + \tau \frac{[\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})][\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})]^T}{(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))^2}$
$\partial_{\mathbf{x}} \partial_{\mathbf{y}} g(\mathbf{x}, \mathbf{y})$	$\partial_{\mathbf{x}} \partial_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}) + \tau \frac{\partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})} + \tau \frac{\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) [\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) - \partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x})]^T}{(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))^2}$
$\frac{\partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x})}{\partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x})}$	$\frac{\partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) - \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) (\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})))^{-1} \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))}{\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))}$

Proposition 2.1. Suppose $\mu, \tau > 0$. If $\mathbf{z}^*(\mathbf{x})$ is a differentiable function in a neighborhood of \mathbf{x} , then $g(\mathbf{x}, \mathbf{y})$ is differentiable and

$$\partial g(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \partial_{\mathbf{x}} g(\mathbf{x}, \mathbf{y}) \\ \partial_{\mathbf{y}} g(\mathbf{x}, \mathbf{y}) \end{pmatrix},$$

where

$$\partial_{\mathbf{x}} g(\mathbf{x}, \mathbf{y}) = \partial_{\mathbf{x}} F(\mathbf{x}, \mathbf{y}) - \tau \frac{\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) - \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})} - \tau \sum_{i=1}^{\kappa} \frac{\partial_{\mathbf{x}} c_i(\mathbf{x})}{c_i(\mathbf{x})}$$

and

$$\partial_{\mathbf{y}} g(\mathbf{x}, \mathbf{y}) = \partial_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}) + \tau \frac{\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})}.$$

Proposition 2.2. Suppose $\mu, \tau > 0$. If $\mathbf{z}^*(\mathbf{x})$ is a differentiable function in a neighborhood of \mathbf{x} and $\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))$ is invertible, then $g(\mathbf{x}, \mathbf{y})$ is twice differentiable and

$$\partial^2 g(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \partial_{\mathbf{x}\mathbf{x}}^2 g(\mathbf{x}, \mathbf{y}) & \partial_{\mathbf{y}} \partial_{\mathbf{x}} g(\mathbf{x}, \mathbf{y}) \\ \partial_{\mathbf{x}} \partial_{\mathbf{y}} g(\mathbf{x}, \mathbf{y}) & \partial_{\mathbf{y}\mathbf{y}}^2 g(\mathbf{x}, \mathbf{y}) \end{pmatrix},$$

where the details of $\partial_{\mathbf{x}\mathbf{x}}^2 g, \partial_{\mathbf{x}} \partial_{\mathbf{y}} g, \partial_{\mathbf{y}\mathbf{y}}^2 g$ are listed in the Table 1.²

We summarise our method in Algorithm 1, which using barrier methods sequentially solves

$$\min_{\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n} F(\mathbf{x}, \mathbf{y}) \quad \text{s.t.} \quad f_{\mu_k}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}) \geq 0, \quad c_i(\mathbf{x}) \geq 0, i \in [\kappa]. \quad (6)$$

In each iteration, for $\tau = \tau_{k,l}$ and $\mu = \mu_k$, we parameterize g by k and l as follows

$$g_{k,l}(\mathbf{x}, \mathbf{y}) \triangleq F(\mathbf{x}, \mathbf{y}) - \tau_{k,l} \ln(f_{\mu_k}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})) - \tau_{k,l} \sum_{i=1}^{\kappa} \ln(c_i(\mathbf{x})).$$

Thus each iteration of our algorithm minimizes the log-barrier function

$$\min_{\mathbf{x}, \mathbf{y}} g_{k,l}(\mathbf{x}, \mathbf{y}). \quad (7)$$

Now we apply the trust-region algorithm to solve problem (7). The trust-region subproblem at the t th step, denoted by $\mathbf{s}_t = (\mathbf{x}_t, \mathbf{y}_t)$, is as follows

$$\min_{\|\mathbf{d}\| \leq \Delta_t} m_{k,l}(\mathbf{d}) = g_{k,l}(\mathbf{s}_t) + \nabla g_{k,l}(\mathbf{s}_t)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T B_t \mathbf{d}. \quad (8)$$

Here B_t needs not to be the exact Hessian of $g_{k,l}$, and the convergence of the trust-region method can still be guaranteed [4]. In practice, we use the Steihaug-Toint truncated CG method [23, 24] to solve

²As the Hessian is symmetric and F and f are twice continuously differentiable, we always have $\partial_{\mathbf{x}} \partial_{\mathbf{y}} g = \partial_{\mathbf{y}} \partial_{\mathbf{x}} g^T$.

Algorithm 1 Bilevel Trust-region Interior-point Method

```

1: Input: parameters  $\eta > 1$ ,  $\tau = \tau_0$ ,  $\mu_0 > 0$ , and  $u > 1$ 
2: for  $k = 1$  to  $K$  do
3:    $\mu = \mu_0/u^k$ ;
4:   for  $l = 1$  to  $T$  do
5:      $\tau = \tau/\eta$ ;
6:     compute  $\mathbf{z}^*(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ ;
7:     derive (approximate) Hessian  $B_k$  for  $g_{k,l}(\mathbf{x}, \mathbf{y})$ ;
8:     solve log-barrier minimization problem (7) to obtain  $(\mathbf{x}_{k,l}, \mathbf{y}_{k,l})$ ;
9:   end for
10: end for

```

the trust-region subproblem, where each step only involves Hessian vector products. This allows us to solve large scale machine learning applications. More details on the complexity analysis for the subproblem see Section 4. We will show in the next section that when μ_k and $\tau_{k,l}$ both approach 0, any limiting point of the sequence generated by Algorithm 1 is a local optimal solution of the original BLO (1) under mild assumptions.

As a closing remark for this section, we point out several main differences of our method with the BVFIM in Liu et al. [13]. First, we simultaneously update \mathbf{x} and \mathbf{y} , while the BVFIM only updates \mathbf{x} . Second, our method is a second-order method, while the BVFIM is a first-order method. Third, only one minimization for $\min_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y})$ is required in each iteration to compute the gradient and Hessian vector products at (\mathbf{x}, \mathbf{y}) in BTRIPM, while the BVFIM needs to solve an additional minimization problem

$$\varphi_{\mu,\tau}(\mathbf{x}) = \min_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}) - \tau \ln(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))$$

to compute the hyper-gradient $\partial_{\mathbf{x}} \phi(\mathbf{x})$. Forth, our method can handle general nonlinear constraint $\mathbf{c}(\mathbf{x}) \geq 0$, while the BVFIM can only handle simply constraints with easy projection if we extend their method using projected gradient methods w.r.t. \mathbf{x} . Moreover, as we will discuss in the next section, we prove that our algorithm converges to a local minimum if each inner sequence (iterates generated by the trust-region method for solving the log-barrier problem (7)) converges to a *local* minimum of each subproblem, while the BVFIM is only guaranteed to converge under a very strong condition that the inner sequence converges to a *global* minimum of each subproblem. Our proof technique is totally different from [13]

3 Theoretical Investigations

In this section, we mainly give theoretical convergence results of the BTRIPM. Before that, let us recall some notation from nonlinear optimization ([18]). Consider a general constrained optimization problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \theta(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{a}(\mathbf{x}) \geq \mathbf{0}, \mathbf{b}(\mathbf{x}) = \mathbf{0}, \end{aligned} \tag{9}$$

where $\theta : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{a} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ and $\mathbf{b} : \mathbb{R}^n \rightarrow \mathbb{R}^l$ are twice continuously differentiable functions in its domain. Denote the Lagrange multipliers of $\mathbf{a}(\mathbf{x}) \geq \mathbf{0}$ and $\mathbf{b}(\mathbf{x}) = \mathbf{0}$ by $\boldsymbol{\lambda}^a$ and $\boldsymbol{\lambda}^b$, respectively. We say $(\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ is an acceptable Lagrange multiplier vector at $\bar{\mathbf{x}}$ if $(\bar{\mathbf{x}}; \boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ satisfies the KKT conditions. Let \mathcal{I} be the active set of the inequality constraints, i.e., $\mathcal{I} \triangleq \{i : \mathbf{a}_i(\bar{\mathbf{x}}) = 0, i \in [p]\}$. We say that the linear independence constraint qualification (LICQ) holds at the feasible point $\bar{\mathbf{x}}$ if $\{\partial a_i(\bar{\mathbf{x}}), \partial b_j(\bar{\mathbf{x}}) : i \in \mathcal{I}, j \in [l]\}$ are linear independent. For problem (9), we say strict complementarity holds at the KKT point $\bar{\mathbf{x}}$ if there exists an acceptable Lagrange multiplier vector $(\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ such that $\bar{\lambda}_i^a > 0$ for all $i \in \mathcal{I}$. Let $L(\mathbf{x}; \boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ denote the Lagrangian function of (9). Define the critical cone

$$\mathcal{C}(\mathbf{x}, \boldsymbol{\lambda}^a) \triangleq \left\{ \mathbf{s} : \begin{array}{l} \partial \mathbf{b}(\mathbf{x}) \mathbf{s} = \mathbf{0}, \partial a_i(\mathbf{x})^T \mathbf{s} = 0 \text{ for all } i \in \mathcal{I} \text{ with } \lambda_i^a > 0, \\ \text{and } \partial a_i(\mathbf{x})^T \mathbf{s} \geq 0 \text{ for all } i \in \mathcal{I} \text{ with } \lambda_i^a = 0 \end{array} \right\}.$$

Let $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ and $(\mathbf{x}_{k,l}^*, \mathbf{y}_{k,l}^*)$ be local minimizers of (6) and (7), respectively. Define $c_0(\mathbf{x}, \mathbf{y}) \triangleq f(\mathbf{x}, \mathbf{y}) - f^*(\mathbf{x})$. We make the following assumption before presenting our main results.

Assumption 3.1. Assume the following conditions hold. (1) The sequence $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ is bounded and, by passing to a subsequence if necessary, $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \rightarrow (\bar{\mathbf{x}}, \bar{\mathbf{y}})$ as $k \rightarrow \infty$. (2) $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \succ 0$. (3) $\mathbf{z}^*(\bar{\mathbf{x}})$ is continuous in a neighborhood of $\bar{\mathbf{x}}$ with $\mathbf{z}^*(\bar{\mathbf{x}}) = \bar{\mathbf{y}}$. (4) The vectors $\{\partial_{\mathbf{x}} c_i(\bar{\mathbf{x}})\}_{i \in \bar{\mathcal{A}}}$ are linearly independent, where $\bar{\mathcal{A}} = \{i : c_i(\bar{\mathbf{x}}) = 0, i = 1, 2, \dots, n\}$

Similar assumptions except (3) are widely used in the IPM literature [6]. Indeed, (3) holds if f is level-bounded in \mathbf{y} , locally uniformly in $\mathbf{x} \in \mathcal{X}$ and $\text{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is a singleton for all $\mathbf{x} \in \mathcal{X}$ (see Lemma 3 in Liu et al. [11]). More discussion on this assumption see Appendix G. We also point out that (4) usually holds in machine learning scenarios, e.g., when $\mathcal{X} = \mathbb{R}^m$ or \mathcal{X} is a box constraint, the linear independence of $\{\partial_{\mathbf{x}} c_i(\bar{\mathbf{x}})\}_{i \in \bar{\mathcal{A}}}$ holds trivially.

Note that $f(\mathbf{x}, \mathbf{y}) \leq f^*(\mathbf{x})$ implies $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = 0$ due to the first-order condition, and, moreover, no CQ holds for (2) ([28]), which means that the KKT conditions may not be necessary optimality conditions of (2). This motivates us to introduce the following auxiliary problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n} \quad & F(\mathbf{x}, \mathbf{y}) \\ \text{s.t.} \quad & c_i(\mathbf{x}) \geq 0, i \in [\kappa], \quad \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}. \end{aligned} \quad (10)$$

Note also that the feasible region of (10) is larger than that of (2) and (10) and (2) have the same objective functions. Under Assumption 3.1, we can show that LICQ holds for (10) (see Lemma D.6), and thus the KKT conditions are necessary optimality conditions for (10). Then using the relationship of the KKT conditions of (10) and (6), we show that the KKT conditions of (10) holds at any limiting point of the local minimum sequence $\{(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)\}$ of (6).

Theorem 3.2. *Let Assumption 3.1 hold. Then for problem (6) and sufficiently large k , the KKT conditions hold at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ with an acceptable Lagrange multiplier vector (λ_0^k, λ^k) . Moreover, letting $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ be any limiting point of the sequence $\{(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)\}$, the KKT conditions hold at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ with an acceptable Lagrange multiplier vector $(\gamma, \nu) = (\lim_{k \rightarrow \infty} \lambda^k, \lim_{k \rightarrow \infty} \lambda_0^k(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)))$ for problem (10), where γ and ν are the Lagrange multipliers of $\mathbf{c}(\mathbf{x}) \geq 0$ and $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = 0$, respectively.*

We further obtain convergence rate results for the BTRIPM in the following theorem.

Theorem 3.3. *Let the conditions of Theorem 3.2 hold and $\bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \gamma, \nu)$ be the Lagrangian function for (10), then there exist Lagrangian multipliers γ_k and ν_k such that*

$$\|\partial \bar{L}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \gamma_k, \nu_k)\| = O(u^{-\frac{k}{2}}). \quad (11)$$

Recalling the definitions of the second-order necessary conditions (SONCs) and second-order sufficient conditions (SOSCs) in optimization literature, which are also stated in Definitions A.4 and A.5 for completeness, we next show the relationship between the KKT solutions of (10) and local minimizers of (1).

Theorem 3.4. *Assume that Assumption 3.1 holds, strict complementarity holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ for (10), and f is third-order continuously differentiable in a neighborhood of $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$. Let $\bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \gamma, \nu)$ be the Lagrangian function for (10) with the same (γ, ν) in Theorem 3.2. Then we have the following conclusions.*

1. *The SONCs hold at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ for problem (10). That is, $\forall \mathbf{s} \in \mathcal{C}((\bar{\mathbf{x}}, \bar{\mathbf{y}}), \gamma)$, we have*

$$\mathbf{s}^T \bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \gamma, \nu) \mathbf{s} \geq \liminf_{k \rightarrow \infty} w_k \|\mathbf{s}\|^2 \geq 0,$$

for some sequence $\{w_k\}$ satisfying $w_k \geq 0 \forall k > 0$, where $\mathcal{C}((\bar{\mathbf{x}}, \bar{\mathbf{y}}), \gamma)$ denotes the critical cone at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$.

2. *If $\limsup_{k \rightarrow \infty} w_k > 0$, then the SOSCs hold at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ for problem (10) and $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ is a strict local optimal solution for (1).*

Following the conditions in Theorem 3.4, we further show the convergence and rate of convergence of BTRIPM to a local optimal solution under suitable choice of T_k .

Theorem 3.5. *Suppose that the assumptions in Theorem 3.4 hold, $\limsup_{k \rightarrow \infty} w_k > 0$ and strict complementarity holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for (6) for sufficiently large k . Then if each T_k is chosen sufficiently large, there exists a subsequence of $\{(\mathbf{x}_{k,l}^*, \mathbf{y}_{k,l}^*)\}$, denoted by $\{(\mathbf{x}_{k',l}^*, \mathbf{y}_{k',l}^*)\}$, such that*

Table 2: Complexity analysis of different algorithms for BLO. We use $\mathbf{p} \in \mathbb{R}^n, \mathbf{q} \in \mathbb{R}^m$ and $\mathbf{d} \in \mathbb{R}^{m+n}$ to denote the intermediate vectors, and Z to denote the intermediate matrix. $J(S)$ denotes the number of CG (truncated CG) steps performed in one outer iteration by CG method (BTRIPM). T_z represents the number of gradient steps to update \mathbf{y} in BVFIM and our method. The analysis for RHG, FHG, CG and BVFIM can be found in [13].

Method	Main Update Steps	Time	Space
RHG	$Z_G^\top \frac{\partial F(\mathbf{x}, \mathbf{y}_G)}{\partial \mathbf{y}}$ with $Z_t = \frac{\partial^2 f}{\partial \mathbf{y}^2} Z_{t-1} + \frac{\partial^2 f}{\partial \mathbf{y} \partial \mathbf{x}}$	$O(cG)$	$O(nG)$
FHG	\mathbf{q}_{-1} with $\mathbf{q}_{t-1} = \mathbf{q}_t + \left(\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{y}} \right)^\top \mathbf{p}_t, \mathbf{p}_{t-1} = \left(\frac{\partial^2 f}{\partial \mathbf{y}^2} \right)^\top \mathbf{p}_t$	$O(cmG)$	$O(mn)$
CG	$-\left(\frac{\partial^2 f(\mathbf{x}, \mathbf{y}_G)}{\partial \mathbf{y} \partial \mathbf{x}} \right)^\top \mathbf{q}$ with $\frac{\partial^2 f}{\partial \mathbf{y}^2} \mathbf{q} = \frac{\partial F}{\partial \mathbf{y}}$ by CG	$O(c(G + J))$	$O(m + n)$
BVFIM	$\frac{\tau_k}{f_{k,l}^{\mathbf{z}} - f(\mathbf{x}_l, \mathbf{y}_{k,l}^{\mathbf{y}})} \left(\frac{\partial f(\mathbf{x}_l, \mathbf{y}_{k,l}^{\mathbf{y}})}{\partial \mathbf{x}} - \frac{\partial f(\mathbf{x}_l, \mathbf{z}_{k,l}^{\mathbf{z}})}{\partial \mathbf{x}} \right)$	$O(c(T_z + T_y))$	$O(m + n)$
Ours	$\min_{\ \mathbf{d}\ \leq \Delta_t} \frac{1}{2} \mathbf{d}^\top B_t \mathbf{d} + \partial g_{k,l}(\mathbf{x}_t, \mathbf{y}_t)^\top \mathbf{d}$ by ST CG	$O(c(T_z + S))$	$O(m + n)$

$\lim_{k' \rightarrow \infty} (\mathbf{x}_{k', T_{k'}}, \mathbf{y}_{k', T_{k'}})$ exists and is a strict local optimal solution of (1). Furthermore, there exist Lagrangian multipliers $\gamma_{k'}$ and $\nu_{k'}$ such that

$$\|\bar{L}(\mathbf{x}_{k', T_{k'}}, \mathbf{y}_{k', T_{k'}}; \gamma_{k'}, \nu_{k'})\| = O(\eta^{-\frac{k'}{2}}) + O(u^{-\frac{k'}{2}}).$$

4 Complexity Analysis

In this section, we analyze the main computational cost of our algorithm and provide a comparison with existing ones. This is based on the following widely used assumptions. (i) The existing methods search the optimal solution of the LL problems by G -step gradient descent. Gradient descent also serves as the transition function of RHG and FHG [13]. (ii) The function or gradient evaluation of F or f , and Hessian-vector product associated with $\frac{\partial^2 f}{\partial \mathbf{y}^2}$ and Jacobian vector product associated with $\frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{y}}$ are supposed to be calculated in time $c = c(m, n)$ for all $\mathbf{p} \in \mathbb{R}^n$ [20, 13].

Now let us consider the complexity of solving our trust-region model (8) using truncated CG. From Table 1, we can see that calculating the exact Hessian vector product associated with $\partial^2 g$ involves solving a linear system³, which is computationally expensive when the dimension of \mathbf{y} is high. This inspires us to use an approximation $B_k(\mathbf{s}_t)$ of the true Hessian, which can still guarantee the convergence of the trust-region method [4]. Specifically, we let $B_k(\mathbf{s}_t) = \partial^2 F(\mathbf{x}_t, \mathbf{y}_t)$, then the computation of Hessian vector product $B_k \mathbf{p}$ is $\partial^2 F(\mathbf{x}_t, \mathbf{y}_t) \mathbf{p}$ for any vector \mathbf{p} . By our second assumption in this section, its computation is in the same order with gradient or function evaluation of F and f .

The main cost in solving the subproblem is the gradient descent steps for $\min_{\mathbf{y}} f(\mathbf{x}_t, \mathbf{y})$ and the Hessian vector products $B_k(\mathbf{s}_t) \mathbf{p}$ in CG steps. Denote by S the number of CG steps taken to entirely solve a trust-region subproblem, and by T_z the number of steps to update \mathbf{y} in $\min_{\mathbf{y}} f(\mathbf{x}_t, \mathbf{y})$. Then the cost of one outer iteration of our algorithm is $O(c(S + T_z))$. Since we only need to restore $\mathbf{z}^*(\mathbf{x}_t)$, \mathbf{s}_t , $B_k(\mathbf{s}_t) \mathbf{p}$ and the gradients, the consumed space is $O(m + n)$. We point out that $B_k(\mathbf{s}_t) \mathbf{p}$ can be directly computed using the structure of BLO in machine learning applications without explicitly formulating the (approximate) Hessian [14]. We summarise the comparison of our algorithm with the existing methods in Table 2.

Although the cost in a single outer iteration of BTRIPM may not be less than the existing algorithms, BTRIPM needs far fewer outer iterations than any existing algorithms. As a second-order algorithm, at most tens of outer iterations are needed for BTRIPM to obtain a highly accurate solution (see next section). In contrast, the existing first-order algorithms usually demand hundreds of outer iterations

³The linear system comes from the Hessian vector product associated with $\partial_{\mathbf{x}\mathbf{x}}^2 f_\mu(\mathbf{x})$. More specifically, for any \mathbf{w} , computing $[\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))]^{-1} [\partial_{\mathbf{x}\mathbf{y}} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))] \mathbf{w}$ involves solving a linear system $[\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))]^{-1} \mathbf{u}$ for $\mathbf{u} = [\partial_{\mathbf{x}\mathbf{y}} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))] \mathbf{w}$.

to reach a relative low accuracy. Therefore, in many cases of interest, BTRIPM is more efficient than the existing methods, as will be illustrated in the next section.

5 Experimental Results

In this section, we conduct experiments to compare the BTRIPM with existing algorithms for BLO. All experiments are implemented using MATLAB R2021b on a PC running Windows 10 Intel(R) Xeon(R) E5-2650 v4 CPU (2.2GHz) and 64GB RAM. Please refer to Appendix E for more experimental results.

5.1 Toy Example

To illustrate the superiority of our algorithm, we test our method on a toy example from Liu et al. [13]

$$\min_{x \in \mathbb{R}, y \in \mathbb{R}} x^2 + y^2, \quad \text{s.t. } y \in \operatorname{argmin}_{y \in \mathbb{R}} \sin(x + y). \quad (12)$$

We underline that the solution set of the LL problem is not a singleton, which prevents most of the existing methods finding the global optimal solution precisely when the initial point is not close to it. We use gradient descent to identify the optimal solution of LL objective $f(x, y) = \sin(x + y)$ in both of BVFIM and our algorithm. Since this problem is simple, we use the exact Hessian as B_k in the trust-region subproblem.

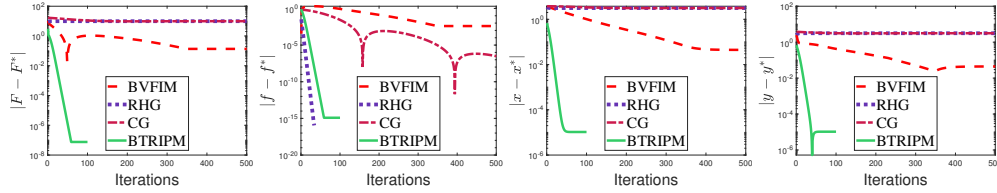


Figure 1: Comparison of BVFIM, RHG, CG with BTRIPM in the toy example (12) with the initial point (x_0, y_0) as $(3, 3)$.

We report comparison results in Figure 1. It shows that BTRIPM enjoys high precision even though the LL problem is nonconvex. RHG and CG fails to find the optimal point though they solve the LL problem well and BVFIM suffers from low accuracy. Furthermore, BTRIPM only requires tens of outer iterations to converge while the existing first-order algorithms needs hundreds.

5.2 Data Hyper-cleaning

To further demonstrate the accuracy and efficiency of our algorithm, we test experiments on MNIST and FasionMNIST [26]. Following the previous experiments in Franceschi et al. [7], we randomly choose 5000 images as the training set, 5000 images as the validation set and 10000 images as the test set. We corrupt 2500 labels in the training set by replacing them with different labels randomly. For the data hyper-cleaning model, we use a single fully-connected layer as the network. Intrinsically, it is a linear model with a linear classifier $\mathbf{y} \in \mathbb{R}^{10} \times \mathbb{R}^{7850}$. After applying a softmax regression, we define the weighted cross entropy loss of the output vector as the LL objective function as the experiment in Liu et al. [13], i.e.,

$$f(\mathbf{x}, \mathbf{y}) = \sum_{(\mathbf{u}_i, v_i) \in \mathcal{D}_{\text{tr}}} \sigma(x_i) \text{CE}(\mathbf{y}, \mathbf{u}_i, v_i).$$

Here (\mathbf{u}_i, v_i) are the sample data, and the input of the sigmoid functions $\mathbf{x} \in \mathbb{R}^{|\mathcal{D}_{\text{tr}}|}$ is regarded as the UL variable, which directly determines the weights. For the UL objective, the original cross entropy loss is used, only depending on the classifier \mathbf{y} , i.e.,

$$F(\mathbf{x}, \mathbf{y}) = \sum_{(\mathbf{u}_i, v_i) \in \mathcal{D}_{\text{val}}} \text{CE}(\mathbf{y}, \mathbf{u}_i, v_i).$$

In data hyper-cleaning model, the weights of the corrupted samples tend to become lower in the training process. Then the corrupted ones will be selected out. In our experiment we choose 0 as the threshold of x_i and calculate the prevalent F1 score of each method. Moreover, we predict the labels

of the validation set by the maximal index of the output vector, thus obtaining the accuracy rate of each method. Then we compare our algorithm with the existing methods using these indexes.

Since BVFIM is obviously stronger than EGBMs and IGBMS in hyper-parameter optimization experiments (see [13] and also evidences from previous subsection), we only compare with BVFIM. We test the two algorithms with two different update rules for parameter $(\mu_k, \tau_{k,l})$. We always set $T = 1$ for the inner iterations and then $\tau_{k,1} = \tau_0/\eta^k$. For BVFIM1 and BTRIPM1, we set $\mu_k = \mu_0/u^k$. For BVFIM2 and BTRIPM2, we update parameters $\mu_k = f(\mathbf{x}_k, \mathbf{y}_k)$ as in [13].

Table 3 shows the accuracy, F1 scores and total wall-clock time of the two methods on two different datasets, where accuracy denotes proportion of labels predicted correctly in the test set and F1 score denotes the harmonic mean of the precision and recall. It can be seen that BTRIPM achieves the most competitive results on two datasets. Furthermore, BTRIPM is faster than BVFIM in both μ updating rules, which indicates robustness of BTRIPM in parameter updating.

Table 3: Comparison of the results of BTRIPM and BVFIM.

Method	MNIST			FashionMNIST		
	Acc.	F1 score	Time(s)	Acc.	F1 score	Time(s)
BVFIM1	0.8829	0.8438	354.2613	0.8402	0.8651	356.1338
BVFIM2	0.9046	0.9186	359.7937	0.8452	0.8923	358.3472
BTRIPM1	0.8981	0.9520	243.5601	0.8428	0.9179	268.8219
BTRIPM2	0.9055	0.9487	241.5207	0.8484	0.9080	262.8357

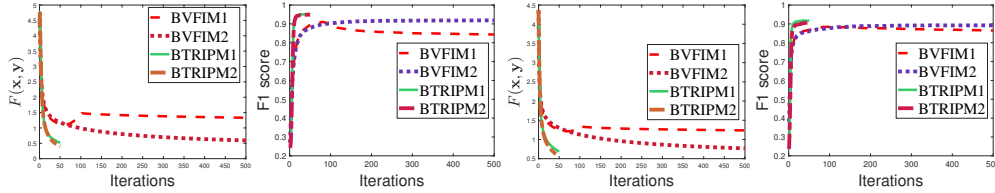


Figure 2: Comparison of the results of BTRIPM and BVFIM for solving data hyper-cleaning tasks. The left two are on MNIST and the right two are on FashionMNIST.

For the same experiment, we also report the UL objective, accuracy and F1 scores of BTRIPM and BVFIM on two datasets in Figure 5.2. As in the last two experiments, BTRIPM converges in tens of outer iterations in both parameter update rules. Meanwhile, BTRIPM achieves higher accuracy and F1 scores.

6 Conclusion and Discussion

In this paper, we propose the first value function based second-order interior-point algorithm for BLO, BTRIPM. We give comprehensive convergence analysis under suitable assumptions that are widely used in the IPM and BLO literature. Computational analysis and numerical experiments have shown the applicability and superiority of our algorithm over existing first-order algorithms. As future work, we would like to apply our algorithm for more practical BLO applications to see if we can use the structure of different models to reduce the truncated CG cost in the trust-region method for solving log-barrier problems.

References

- [1] Ben-Ayed, O. and Blair, C. E. Computational difficulties of bilevel linear programming. *Operations Research*, 38(3):556–560, 1990.
- [2] Bishop, N., Tran-Thanh, L., and Gerding, E. Optimal learning from verified training data. In *Advances in Neural Information Processing Systems*, volume 33, pp. 9520–9529, 2020.
- [3] Brückner, M. and Scheffer, T. Stackelberg games for adversarial prediction problems. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 547–555, 2011.

- [4] Conn, A. R., Gould, N. I., and Toint, P. L. *Trust region methods*. SIAM, 2000.
- [5] Dempe, S. Bilevel optimization: theory, algorithms, applications and a bibliography. In *Bilevel Optimization*, pp. 581–672. Springer, 2020.
- [6] Forsgren, A., Gill, P. E., and Wright, M. H. Interior methods for nonlinear optimization. *SIAM Review*, 44(4):525–597, 2002.
- [7] Franceschi, L., Donini, M., Frasconi, P., and Pontil, M. Forward and reverse gradient-based hyperparameter optimization. In *International Conference on Machine Learning*, pp. 1165–1173. PMLR, 2017.
- [8] Franceschi, L., Frasconi, P., Salzo, S., Grazzi, R., and Pontil, M. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pp. 1568–1577. PMLR, 2018.
- [9] Ji, K., Yang, J., and Liang, Y. Bilevel optimization: Convergence analysis and enhanced design. In *International Conference on Machine Learning*, pp. 4882–4892. PMLR, 2021.
- [10] Johnson, B. A. and Iizuka, K. Integrating openstreetmap crowdsourced data and landsat time-series imagery for rapid land use/land cover (lulc) mapping: Case study of the laguna de bay area of the philippines. *Applied Geography*, 67:140–149, 2016.
- [11] Liu, R., Mu, P., Yuan, X., Zeng, S., and Zhang, J. A generic first-order algorithmic framework for bi-level programming beyond lower-level singleton. In *International Conference on Machine Learning*, pp. 6305–6315. PMLR, 2020.
- [12] Liu, R., Gao, J., Zhang, J., Meng, D., and Lin, Z. Investigating bi-level optimization for learning and vision from a unified perspective: A survey and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021 doi: 10.1109/TPAMI.2021.3132674.
- [13] Liu, R., Liu, X., Yuan, X., Zeng, S., and Zhang, J. A value-function-based interior-point method for non-convex bi-level optimization. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 6882–6892. PMLR, 2021.
- [14] Lorraine, J., Vicol, P., and Duvenaud, D. Optimizing millions of hyperparameters by implicit differentiation. In *International Conference on Artificial Intelligence and Statistics*, pp. 1540–1552. PMLR, 2020.
- [15] MacKay, M., Vicol, P., Lorraine, J., Duvenaud, D., and Grosse, R. Self-tuning networks: Bilevel optimization of hyperparameters using structured best-response functions. In *International Conference on Learning Representations (ICLR)*, 2019.
- [16] Magnus, J. R. and Neudecker, H. *Matrix differential calculus with applications in statistics and econometrics*. John Wiley & Sons, 2019.
- [17] Mehrlitz, P. and Zemkoho, A. B. Sufficient optimality conditions in bilevel programming. *Mathematics of Operations Research*, 46(4):1573–1598, 2021.
- [18] Nocedal, J. and Wright, S. *Numerical optimization*. Springer Science & Business Media, 2006.
- [19] Pedregosa, F. Hyperparameter optimization with approximate gradient. In *International Conference on Machine Learning*, pp. 737–746. PMLR, 2016.
- [20] Rajeswaran, A., Finn, C., Kakade, S. M., and Levine, S. Meta-learning with implicit gradients. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- [21] Rohra, J. G., Perumal, B., Narayanan, S. J., Thakur, P., and Bhatt, R. B. User localization in an indoor environment using fuzzy hybrid of particle swarm optimization & gravitational search algorithm with neural networks. In *Proceedings of Sixth International Conference on Soft Computing for Problem Solving*, pp. 286–295. Springer, 2017.
- [22] Shaban, A., Cheng, C.-A., Hatch, N., and Boots, B. Truncated back-propagation for bilevel optimization. In *AISTATS*, 2019.

- 327 [23] Steihaug, T. The conjugate gradient method and trust regions in large scale optimization. *SIAM*
328 *Journal on Numerical Analysis*, 20(3):626–637, 1983.
- 329 [24] Toint, P. Towards an efficient sparsity exploiting newton method for minimization. In *Sparse*
330 *matrices and their uses*, pp. 57–88. Academic press, 1981.
- 331 [25] Wang, J., Chen, H., Jiang, R., Li, X., and Li, Z. Fast algorithms for stackelberg prediction
332 game with least squares loss. In *Proceedings of the 38th International Conference on Machine*
333 *Learning*, pp. 10708–10716. PMLR.
- 334 [26] Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking
335 machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 336 [27] Yang, Z., Chen, Y., Hong, M., and Wang, Z. Provably global convergence of actor-critic: A
337 case for linear quadratic regulator with ergodic cost. In *Proceedings of the 33rd International*
338 *Conference on Neural Information Processing Systems*, pp. 8353–8365, 2019.
- 339 [28] Ye, J. J., Yuan, X., Zeng, S., and Zhang, J. Difference of convex algorithms for bilevel programs
340 with applications in hyperparameter selection. *arXiv preprint arXiv:2102.09006*, 2021.
- 341 [29] Ye, J. J. and Zhu, D. Optimality conditions for bilevel programming problems. *Optimization*,
342 33:9–27, 1995.

343 Checklist

- 344 1. For all authors...
- 345 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
346 contributions and scope? [\[Yes\]](#) Please refer to Section 1 for detailed description
347 corresponding to the abstract.
- 348 (b) Did you describe the limitations of your work? [\[No\]](#)
- 349 (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
- 350 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
351 them? [\[Yes\]](#)
- 352 2. If you are including theoretical results...
- 353 (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) Please refer to
354 Section 3 in the main body and Section ?? in the supplementary materials.
- 355 (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) We put most of the
356 detailed proofs in the supplementary materials.
- 357 3. If you ran experiments...
- 358 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
359 mental results (either in the supplemental material or as a URL)? [\[Yes\]](#) Instruction for
360 the code, datasets are provided in Section 5 and the supplemental material.
- 361 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
362 were chosen)? [\[Yes\]](#) Information about the dataset and hyperparameters for numerical
363 experiments can be found in Section 5 and the supplemental material.
- 364 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
365 ments multiple times)? [\[N/A\]](#) Some experiments e.g. problem 12 is deterministic.
- 366 (d) Did you include the total amount of compute and the type of resources used (e.g., type
367 of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) Please refer to Section 5 for detail
368 information.
- 369 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 370 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
- 371 (b) Did you mention the license of the assets? [\[Yes\]](#)
- 372 (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#)
- 373 (d) Did you discuss whether and how consent was obtained from people whose data you’re
374 using/curating? [\[N/A\]](#)

- 375 (e) Did you discuss whether the data you are using/curating contains personally identifiable
376 information or offensive content? [N/A] All the experiments in this paper are based on
377 public data.
- 378 5. If you used crowdsourcing or conducted research with human subjects...
- 379 (a) Did you include the full text of instructions given to participants and screenshots, if
380 applicable? [N/A]
- 381 (b) Did you describe any potential participant risks, with links to Institutional Review
382 Board (IRB) approvals, if applicable? [N/A]
- 383 (c) Did you include the estimated hourly wage paid to participants and the total amount
384 spent on participant compensation? [N/A]

Appendix

The appendix is organized as follows. We firstly introduce some necessary definitions and theorems in nonlinear optimization for our proofs in Appendix A. In Appendices B to D, we prove the theoretical conclusions in Sections 2 and 3. We validate all of the assumptions in this paper in Appendix G. For numerical experiments, we provide supplementary experiments in Appendix E and the details of all experiments in Appendix F.

A Definitions and Technical Lemmas from Nonlinear Programming

We supplement some basic definitions and theorems in nonlinear optimization for problem (9) from Forsgren et al. [6], Nocedal & Wright [18].

Definition A.1. We say the Mangasarian–Fromovitz constraint qualification (MFCQ) holds at the feasible point $\bar{\mathbf{x}}$ if $\bar{\mathbf{x}}$ if there exists a vector \mathbf{p} such that $\partial a_i(\bar{\mathbf{x}})^T \mathbf{p} > 0$ for all $i \in \mathcal{I}$ and $\partial \mathbf{b}(\bar{\mathbf{x}})^T \mathbf{p} = 0$ for all $i \in \mathcal{I}$.

Theorem A.2. (First order Necessary Conditions) Suppose $\bar{\mathbf{x}}$ is a local optimal point of (9) where LICQ holds. Then there is a Lagrange multiplier vector $(\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ such that the following conditions hold

$$\begin{aligned} \partial \theta(\bar{\mathbf{x}}) - \sum_{i \in \mathcal{I}} \lambda_i^a \partial a_i(\bar{\mathbf{x}}) - \sum_{i=1}^l \lambda_i^b \partial b_i(\bar{\mathbf{x}}) &= \mathbf{0}, \\ \mathbf{a}(\bar{\mathbf{x}}) &\geq \mathbf{0}, \\ \mathbf{b}(\bar{\mathbf{x}}) &= \mathbf{0}, \\ \boldsymbol{\lambda}^a &\geq \mathbf{0}, \\ \lambda_i^a a_i(\bar{\mathbf{x}}) &= 0 \text{ for all } i \in [p]. \end{aligned} \tag{13}$$

The conditions (13) are known as the KKT conditions.

Definition A.3. For problem (9), we say strict complementarity holds at the KKT point $\bar{\mathbf{x}}$ if there exists an acceptable Lagrange multiplier vector $(\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ such that $\bar{\lambda}_i^a > 0$ for all $i \in \mathcal{I}$.

Definition A.4. We say $\bar{\mathbf{x}}$ satisfies the second-order necessary conditions (SONCs) if $\bar{\mathbf{x}}$ is a KKT point and there exists an acceptable Lagrange multiplier $(\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ such that for all $\mathbf{s} \in \mathcal{C}(\bar{\mathbf{x}}, \boldsymbol{\lambda}^a)$,

$$\mathbf{s}^T \partial^2 L(\bar{\mathbf{x}}; \boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b) \mathbf{s} \geq 0.$$

Definition A.5. We say $\bar{\mathbf{x}}$ satisfies the second-order sufficient conditions (SOSCs) if $\bar{\mathbf{x}}$ is a KKT point and there exists an acceptable Lagrange multiplier $(\boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b)$ such that for all $\mathbf{s} \in \mathcal{C}(\bar{\mathbf{x}}, \boldsymbol{\lambda}^a) \setminus \{\mathbf{0}\}$,

$$\mathbf{s}^T \partial^2 L(\bar{\mathbf{x}}; \boldsymbol{\lambda}^a, \boldsymbol{\lambda}^b) \mathbf{s} > 0.$$

We have the following Theorems.

Theorem A.6 (Theorem 12.5 in Nocedal & Wright [18]). Suppose that $\bar{\mathbf{x}}$ is a local optimal point of (9) and that LICQ holds. Then $\bar{\mathbf{x}}$ satisfies the SONCs.

Theorem A.7 (Theorem 12.6 in Nocedal & Wright [18]). If $\bar{\mathbf{x}}$ satisfies the SOSCs, then $\bar{\mathbf{x}}$ is a strict local optimal point for (9).

Now we consider problem (6). The following Assumption is normal in IPM literature, which can be implied by the conditions of Theorem 3.4 (see Lemma D.9).

Assumption A.8. Assume that (1) the MFCQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$, (2) there exists $w_k > 0$ such that $\mathbf{p}^T \partial^2 L(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \boldsymbol{\lambda}) \mathbf{p} \geq w_k \|\mathbf{p}\|^2$ for all acceptable $\boldsymbol{\lambda}$ and all nonzero \mathbf{p} satisfying $\partial F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^T \mathbf{p} = 0$ and $J_{\mathcal{I}}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{p} \geq 0$, where $L(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \boldsymbol{\lambda})$ denotes the Lagrangian function and \mathcal{I} denotes the active index set of constraints in (6), and (3) strict complementarity hold at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ in (6).

The following two Lemmas are directly from Theorem 3.12, Lemma 3.13 of [6], which will be used in our proofs for Section 3.

Lemma A.9. Assume that a log-barrier method is applied for (6), in which $\tau_{k,l}$ monotonically converges to zero as $l \rightarrow \infty$. Assume that $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ is a local minimizer of (6) and (1) and (2) of Assumption A.8 hold. Then for any $k > 0$, there exists at least one subsequence of unconstrained minimizers of the barrier functions $g_{k,l}(\mathbf{x}, \mathbf{y}) = F(\mathbf{x}, \mathbf{y}) - \tau_{k,l} \ln(f_{\mu_k}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})) - \tau_{k,l} \sum_{i=1}^{\kappa} \ln(c_i(\mathbf{x}))$ converging to $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ as $l \rightarrow \infty$.

425 **Lemma A.10.** Consider problem (6). Assume the conditions of Lemma A.9 are satisfied, and
 426 additionally the strict complementarity holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$. Let $\{(\mathbf{x}_{k,l}^*, \mathbf{y}_{k,l}^*)\}$ denote the converging
 427 subsequence in Lemma A.9, then $\|(\mathbf{x}_{k,l}^*, \mathbf{y}_{k,l}^*) - (\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)\| = O(\tau_{k,l})$.

428 B Proof of Section 2

429 B.1 Proof of Proposition 2.1

430 *Proof.* The only difficulty lies in computing $\partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x})$. According to the differentiability of $\mathbf{z}^*(\mathbf{x})$
 431 and $f(\mathbf{x}, \mathbf{y})$, $f_{\mu}^*(\mathbf{x}) = f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) + \mu$ is also continuously differentiable w.r.t. \mathbf{x} . By the chain rule,
 432 we have

$$\partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) + \left(\frac{\partial \mathbf{z}^*(\mathbf{x})}{\partial \mathbf{x}} \right)^T \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})).$$

433 Since $\mathbf{z}^*(\mathbf{x}) \in \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{x}, \mathbf{y})$, we have $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) = \mathbf{0}$ from the first-order optimality
 434 condition. So we have $\partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))$. \square

435 B.2 Proof of Proposition 2.2

436 *Proof.* The main computational complexity concentrates upon the derivative of $\ln(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))$.
 437 First, we can calculate the first-order derivative applying the conclusion of Proposition 2.1. By the
 438 chain rule, we have

$$\begin{aligned} \partial_{\mathbf{x}} [\ln(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))] &= \frac{\partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x}) - \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})}, \\ \partial_{\mathbf{y}} [\ln(f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}))] &= \frac{-\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})}{f_{\mu}^*(\mathbf{x}) - f(\mathbf{x}, \mathbf{y})}. \end{aligned}$$

439 Then we can calculate the second-order derivatives $\partial_{\mathbf{x}\mathbf{x}}^2 g, \partial_{\mathbf{x}} \partial_{\mathbf{y}} g, \partial_{\mathbf{y}\mathbf{y}}^2 g$ as in Table 1. Note that
 440 $\partial_{\mathbf{x}} \partial_{\mathbf{y}} g = \partial_{\mathbf{y}} \partial_{\mathbf{x}} g^T$ because of second-order differentiability. The only difficulty lies in computing
 441 $\partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x})$.
 442 According to the proof in Proposition 2.1, we have

$$\partial_{\mathbf{x}} f_{\mu}^*(\mathbf{x}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) + \frac{\partial \mathbf{z}^*(\mathbf{x})}{\partial \mathbf{x}}^T \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})).$$

443 Note that we always have $\frac{\partial \mathbf{z}^*(\mathbf{x})}{\partial \mathbf{x}}^T \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) = \mathbf{0}$ holds for all \mathbf{x} as $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) = \mathbf{0}$.
 444 Therefore, we have

$$\partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x}) = \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) + \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) \frac{\partial \mathbf{z}^*(\mathbf{x})}{\partial \mathbf{x}}.$$

445 Applying the implicit function theorem to $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) = \mathbf{0}$ w.r.t. \mathbf{x} , we have

$$\partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) + \partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) \frac{\partial \mathbf{z}^*(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{0}_{n \times m}.$$

446 So we can calculate $\partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x})$ as

$$\partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x}) = \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) - \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) (\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})))^{-1} \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})).$$

447 The remaining calculation is obvious. \square

448 C Proof of Section 3

449 In the proofs of this section, we make the following convention. Let $(\lambda_0^k, \boldsymbol{\lambda}^k)$, where $\boldsymbol{\lambda}^k =$
 450 $(\lambda_1^k, \dots, \lambda_{\kappa}^k)^T$, be an acceptable Lagrange multiplier for problem (6), where λ_0^k associates with
 451 $c_0(\mathbf{x}, \mathbf{y}) \leq \mu_k$, and λ_i^k associates with $c_i(\mathbf{x}) \geq 0$ $i \in [\kappa]$. Now let us define active sets and Jacobian
 452 for $\mathbf{c}(\mathbf{x}) \geq 0$:

$$\begin{aligned} \mathcal{A}_k &\triangleq \{i : c_i(\mathbf{x}_k), i \in [\kappa]\}, \\ \bar{\mathcal{A}} &\triangleq \{i : c_i(\bar{\mathbf{x}}) = 0, i \in [\kappa]\}, \\ J_{\mathcal{E}}(\mathbf{x}) &\triangleq (\partial_{\mathbf{x}} c_{i_1}(\mathbf{x}), \partial_{\mathbf{x}} c_{i_2}(\mathbf{x}), \dots, \partial_{\mathbf{x}} c_{i_t}(\mathbf{x}))^T, i_j \in \mathcal{E}, \end{aligned}$$

453 where $\mathcal{E} \in [\kappa]$ is an arbitrary index set.

454 **Lemma C.1.** *Let the notation be the same with Section 3. Then $\bar{\mathbf{y}} \in \operatorname{argmin}_{\mathbf{y}} f(\bar{\mathbf{x}}, \mathbf{y})$.*

455 *Proof.* Because of the continuity of $\mathbf{z}^*(\mathbf{x})$ in a neighborhood of $\bar{\mathbf{x}}$, $f^*(\mathbf{x}) = f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))$ is a
 456 continuous function w.r.t. \mathbf{x} in this neighborhood. Since $f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - f^*(\bar{\mathbf{x}}_k) \leq \mu_k$ and $\lim_{k \rightarrow \infty} \mu_k = 0$,
 457 we have $f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = f^*(\bar{\mathbf{x}})$. \square

458 C.1 Proof of Theorem 3.2

459 *Proof.* Since $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ is a local optimal solution of problem (6) and, by Lemma D.5, LICQ holds
 460 when k is large enough, we have the following KKT conditions in problem (6) if k is sufficiently
 461 large,

$$\begin{aligned} \partial F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \sum_{i=1}^{\kappa} \lambda_i^k \partial c_i(\bar{\mathbf{x}}_k) + \lambda_0^k [\partial f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial f^*(\bar{\mathbf{x}}_k)] &= \mathbf{0}, \\ c_i(\bar{\mathbf{x}}_k) &\geq 0, i = 1, 2, \dots, \kappa, \\ c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) &\leq \mu_k, \\ \lambda_i^k &\geq 0, i = 0, 1, \dots, \kappa, \\ \lambda_i^k c_i(\bar{\mathbf{x}}_k) &= 0, i = 0, 1, \dots, \kappa. \end{aligned} \quad (14)$$

462 Noting that $\partial_{\mathbf{y}} c_i(\mathbf{x}) = \mathbf{0}$ for $i \in [\kappa]$, $\partial_{\mathbf{y}} f^*(\mathbf{x}) = \mathbf{0}$ and $\partial_{\mathbf{y}} c_0(\mathbf{x}, \mathbf{y}) = \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$, we have

$$\lambda_0^k \partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = -\partial_{\mathbf{y}} F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k). \quad (15)$$

463 In Assumption 3.1 we assume $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \rightarrow (\bar{\mathbf{x}}, \bar{\mathbf{y}})$ as $k \rightarrow \infty$. Due to the definition of $\mathbf{z}^*(\bar{\mathbf{x}}_k)$ and
 464 Lemma C.1, we have

$$\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k) \rightarrow \mathbf{0}. \quad (16)$$

465 By Taylor expansion, we have

$$\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = \partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) + \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) [\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)] + \mathbf{v}_k, \quad (17)$$

466 for some \mathbf{v}_k satisfying $\|\mathbf{v}_k\| = O(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|^2)$. Moreover, since $\mathbf{z}^*(\bar{\mathbf{x}}_k) = \operatorname{argmin}_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y})$,
 467 by first-order condition we have

$$\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) = \mathbf{0}. \quad (18)$$

468 Since $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \succ 0$, let a_1 and a_2 be such that where $0 < a_1 < a_2$, $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \succeq a_1 I$ and
 469 $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \preceq a_2 I$. By $\lim_{k \rightarrow \infty} (\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = (\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we have

$$\frac{a_1}{2} I \preceq \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \preceq 2a_2 I$$

470 when k is large enough. Then we have

$$\frac{a_1}{2} \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| \leq \|\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) [\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)]\| \leq 2a_2 \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|,$$

471 or equivalently, $\|\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))\| = \Theta(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|)$. This, together with (17)
 472 and (18), implies $\|\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)\| = \Theta(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|)$. Therefore, letting $k \rightarrow \infty$ in (15) and using
 473 (16), we deduce $\lambda_0^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| = \Theta(1)$ when $\partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \neq \mathbf{0}$, and $\lambda_0^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| = o(1)$
 474 when $\partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \mathbf{0}$. We first discuss the situation that $\partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \neq \mathbf{0}$. Since $\{\lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))\}$
 475 is bounded, by passing to a subsequence if necessary, we assume that $\{\lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))\}$ converges
 476 to a limiting point (say, \mathbf{q}). By substituting (17) to (15), noting (18) and letting $k \rightarrow \infty$, we have

$$\lambda_0^k \mathbf{v}_k \rightarrow -\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{q} - \partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}). \quad (19)$$

477 Using (16) and the definition of \mathbf{v}_k , we have

$$\|\lambda_0^k \mathbf{v}_k\| = O(\|\lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))\| \cdot \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|) \rightarrow O(\|\mathbf{q}\| \cdot \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|) \rightarrow 0.$$

478 This, together with (19), implies

$$\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{q} = -\partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}).$$

479 On the other hand, it follows from (14) that

$$\partial_{\mathbf{x}} F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \sum_{i=1}^{\kappa} \lambda_i^k \partial_{\mathbf{x}} c_i(\bar{\mathbf{x}}_k) + \lambda_0^k [\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k)] = \mathbf{0}. \quad (20)$$

480 Note that $\partial_{\mathbf{x}} f^*(\mathbf{x}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x}))$, we can also use a similar Taylor expansion on $\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$,

$$\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k) = [\partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k))](\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)) + \mathbf{u}_k \quad (21)$$

481 for some \mathbf{u}_k satisfying $\|\mathbf{u}_k\| = O(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|^2)$. Then (21), together with (16), implies

$$\begin{aligned} & \|\lambda_0^k [\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k)]\| \\ & \leq \|\partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) \lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))\| + \|\mathbf{u}_k\| \\ & = O(\|\partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) \mathbf{q}\|) = O(1). \end{aligned}$$

482 Then (20) implies

$$\left\| \sum_{i=1}^{\kappa} \lambda_i^k \partial_{\mathbf{x}} c_i(\bar{\mathbf{x}}_k) \right\| = \|\partial_{\mathbf{x}} F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) + \lambda_0^k [\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k)]\| \leq \|\partial_{\mathbf{x}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}})\| + O(1) \quad (22)$$

483 is uniformly bounded for any k . Using (21) and (16), we obtain

$$\lim_{k \rightarrow \infty} \lambda_0^k [\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k)] = \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{q}.$$

484 Due to Lemma D.7, $\{\lambda_i^k\} \forall i \in [\kappa]$ is bounded. Therefore for all $i \in [\kappa]$, limiting points of $\{\lambda_i^k\}$ (say, λ_i) exist. Then from (16), (20) and (21), by subsequence convergence we have

$$\partial_{\mathbf{x}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) - \sum_{i=1}^{\kappa} \lambda_i \partial_{\mathbf{x}} c_i(\bar{\mathbf{x}}) + \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{q} = \mathbf{0}. \quad (23)$$

486 Moreover, by Lemma (D.5), we know LICQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for (6), then the acceptable Lagrange multiplier $(\lambda_0^k, \boldsymbol{\lambda}^k)$ at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ must be unique. Similarly, by Lemma D.6 we know LICQ holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ for (10), then the acceptable Lagrange multiplier at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ must be unique. Then the limiting point of $\{(\lambda_0^k, \boldsymbol{\lambda}^k)\}$ is unique. Thus we have $(\boldsymbol{\gamma}, \boldsymbol{\nu}) = (\lim_{k \rightarrow \infty} \boldsymbol{\lambda}^k, \lim_{k \rightarrow \infty} \lambda_0^k(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)))$, where $\boldsymbol{\gamma} = (\lambda_1 \ \cdots \ \lambda_{\kappa})^T$ and $\boldsymbol{\nu} = \mathbf{q}$.

491 By combining (23) and (15), it follows that

$$\partial F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) - \sum_{i=1}^{\kappa} \lambda_i \partial c_i(\bar{\mathbf{x}}) + \begin{pmatrix} \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \\ \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} \mathbf{q} = \mathbf{0}, \quad (24)$$

492 where $\partial c_i(\bar{\mathbf{x}}) = \begin{pmatrix} \partial_{\mathbf{x}} c_i(\bar{\mathbf{x}}) \\ \mathbf{0} \end{pmatrix}$ and $\begin{pmatrix} \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \\ \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} = (\partial(\partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}})))^T$ as f is twice continuously differentiable in a neighborhood of $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$. Note that

$$c_i(\bar{\mathbf{x}}) \geq 0, i \in [\kappa] \quad (25)$$

$$c_0(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \leq 0, \quad (26)$$

$$\lambda_i \geq 0, i \in [\kappa], \quad (27)$$

$$\lambda_i c_i(\bar{\mathbf{x}}) = 0, \quad (28)$$

494 can be directly obtained by deducing the limit of (14). Note that (26) implies

$$\partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \mathbf{0}. \quad (29)$$

495 Then (25), (29), (27), (28) and (24) consist of the KKT conditions at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ for (10).

496 When $\partial_{\mathbf{y}} \bar{F}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \mathbf{0}$, $\lambda_0^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| = o(1)$. Then by using Taylor expansion (17) and (21) and letting $k \rightarrow \infty$ in (14), we have

$$\partial F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) - \sum_{i=1}^{\kappa} \lambda_i \partial c_i(\bar{\mathbf{x}}) = \mathbf{0},$$

498 which is exactly (24) with $\mathbf{q} = \mathbf{0}$. Note that (25), (29), (27) and (28) also hold for this situation. \square

499 **C.2 Proof of Theorem 3.3**

500 We use the same notation with the proof of Theorem 3.2. Substitute (16) and (21), in the first equation
501 of (14), we have

$$\partial F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \sum_{i=1}^{\kappa} \lambda_i^k \partial c_i(\bar{\mathbf{x}}_k) + \lambda_0^k \partial_{\mathbf{y}} \partial f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)(\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)) = O(\lambda_0^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|^2).$$

502 As we have proved $\lambda_0^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| = O(1)$ and $\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| = O(\mu_k^{\frac{1}{2}})$ in the proof of Theorem
503 3.2, we thus obtain that

$$\partial F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \sum_{i=1}^{\kappa} \lambda_i^k \partial c_i(\bar{\mathbf{x}}_k) + \lambda_0^k \partial_{\mathbf{y}} \partial f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)(\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)) = O(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|) = O(\mu_k^{\frac{1}{2}}).$$

504 Let $\boldsymbol{\gamma}_k = (\lambda_1^k, \lambda_2^k, \dots, \lambda_{\kappa}^k)^T$, $\boldsymbol{\nu}_k = \lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))$. Then by the definition of Lagrangian function,
505 we have

$$\partial \bar{L}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \boldsymbol{\gamma}_k, \boldsymbol{\nu}_k) = \partial F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \sum_{i=1}^{\kappa} \lambda_i^k \partial c_i(\bar{\mathbf{x}}_k) + \lambda_0^k \partial_{\mathbf{y}} \partial f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)(\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)).$$

506 Therefore we have $\partial \bar{L}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \boldsymbol{\gamma}_k, \boldsymbol{\nu}_k) = O(\mu_k^{\frac{1}{2}}) = O(u^{-\frac{\kappa}{2}})$, which completes the proof.

507 **C.3 Proof of Theorem 3.4**

508 *Proof.* From Lemma D.5 we know under Assumption 3.1, the LICQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for problem
509 (6) when k is large enough. Theorem 3.2 shows that $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ must be a KKT point for problem (6).
510 Then from Theorem A.6 we obtain that the SONCs holds on $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for problem (6).
511 To show second-order necessary conditions holds, it suffices to show for all \mathbf{s} satisfying

$$\partial c_{i_j}(\bar{\mathbf{x}})^T \mathbf{s} = 0, \quad \forall i_j \in \bar{\mathcal{A}}, \quad (30)$$

$$\begin{pmatrix} \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \\ \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} \mathbf{s} = \mathbf{0}, \quad (31)$$

512 we always have

$$\mathbf{s}^T \partial^2 \bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \boldsymbol{\gamma}, \boldsymbol{\nu}) \mathbf{s} \geq 0, \quad (32)$$

513 where

$$\bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \boldsymbol{\gamma}, \boldsymbol{\nu}) = F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + \mathbf{c}(\bar{\mathbf{x}})^T \boldsymbol{\gamma} + (\partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^T \boldsymbol{\nu} \quad (33)$$

514 is the Lagrangian function. Define

$$P(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \begin{pmatrix} I_m & 0_{m \times n} \\ (\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^{-1} \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & I_n \end{pmatrix},$$

515 and $\mathbf{l} = P(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{s}$. Then (31) is equivalent to

$$\begin{pmatrix} \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \\ \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix}^T P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^{-1} \mathbf{l} = \mathbf{0}.$$

516 This further yields

$$\begin{pmatrix} 0_{m \times n} \\ \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix}^T \mathbf{l} = \mathbf{0}.$$

517 Set $\mathbf{l} = (\mathbf{l}_{\mathbf{x}}^T \quad \mathbf{l}_{\mathbf{y}}^T)^T$. Then by the positive definiteness of $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we conclude (30) is equivalent
518 to

$$\mathbf{l}_{\mathbf{y}} = \mathbf{0}. \quad (34)$$

519 On the other hand, since $\partial c_i(\bar{\mathbf{x}}) = \begin{pmatrix} \partial_{\mathbf{x}} c_i(\bar{\mathbf{x}}) \\ \mathbf{0} \end{pmatrix}$, (30) is equivalent to $(J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) \quad 0_{t \times n}) P^{-1} \mathbf{l} = \mathbf{0}$,

520 where $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) = (\partial_{\mathbf{x}} c_{i_1}(\bar{\mathbf{x}}) \quad \dots \quad \partial_{\mathbf{x}} c_{i_t}(\bar{\mathbf{x}}))^T$. Noting that

$$P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^{-1} = \begin{pmatrix} I_m & 0_{m \times n} \\ -(\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^{-1} \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & I_n \end{pmatrix},$$

521 we have

$$(J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) \quad 0_{t \times n}) P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^{-1} = (J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) \quad 0_{t \times n}),$$

522 and thus (30) is finally equivalent to

$$J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) \mathbf{l}_{\mathbf{x}} = \mathbf{0}. \quad (35)$$

523 In summary, (31) and (30) are equivalent to (34) and (35).

524 Now we consider the inequality (32). Using the expression

$$\partial^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \begin{pmatrix} \partial_{\mathbf{x}\mathbf{x}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \\ \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}})^T & \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix},$$

525 we then have

$$\partial^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^T \begin{pmatrix} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & 0_{m \times n} \\ 0_{n \times m} & \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} P(\bar{\mathbf{x}}, \bar{\mathbf{y}}).$$

526 For \mathbf{s} satisfying (31) and (30), we have

$$\mathbf{s}^T \partial_{y_i} \partial^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{s} = 2\mathbf{s}^T \partial_{y_i} P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^T \begin{pmatrix} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & 0_{m \times n} \\ 0_{n \times m} & \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} \mathbf{1} + \mathbf{1}^T \begin{pmatrix} \partial_{y_i} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & 0_{m \times n} \\ 0_{n \times m} & \partial_{y_i} \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} \mathbf{1}. \quad (36)$$

527 As we have already shown $\mathbf{l}_{\mathbf{y}} = \mathbf{0}$, using $\mathbf{s} = P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^{-1} \mathbf{1}$ and $\partial_{y_i} P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^T = \begin{pmatrix} 0_{m \times m} & * \\ 0_{n \times m} & 0_{n \times n} \end{pmatrix}$,

528 we have $\mathbf{s}^T \partial_{y_i} P(\bar{\mathbf{x}}, \bar{\mathbf{y}})^T \begin{pmatrix} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & 0_{m \times n} \\ 0_{n \times m} & \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \end{pmatrix} \mathbf{1} = 0$ and thus from (36)

$$\mathbf{s}^T \partial_{y_i} \partial^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{s} = \mathbf{l}_{\mathbf{x}}^T \partial_{y_i} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{l}_{\mathbf{x}}. \quad (37)$$

529 By third order contentiously differentiability of f at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we have

$$\partial^2 \partial_{y_i} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \partial_{y_i} \partial^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}). \quad (38)$$

530 By (38), (37), (33), Lemma D.10 and the expression $(\gamma, \nu) =$
 531 $(\lim_{k \rightarrow \infty} \lambda^k, \lim_{k \rightarrow \infty} \lambda_0^k(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)))$, we obtain

$$\mathbf{s}^T \partial^2 \bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \gamma, \nu) \mathbf{s} \geq \liminf_{k \rightarrow \infty} w_k \|\mathbf{s}\|^2 \geq 0. \quad (39)$$

532 We point out that as strict complementarity holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we always have $\gamma_{\bar{\mathcal{A}}} > 0$. Then by
 533 definition, we have the critical cone

$$\mathcal{C}((\bar{\mathbf{x}}, \bar{\mathbf{y}}), \gamma) = \{\mathbf{s} : \partial(\partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}})) \mathbf{s} = \mathbf{0}, \partial c_i(\mathbf{x})^T \mathbf{s} = 0 \text{ for all } i \in \bar{\mathcal{A}}\}.$$

534 This, together with (39) and Definition A.4, implies that the SONCs holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$.

535 For the second part of the conclusions, note that when strict complementarity holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ we only
 536 need to prove

$$\mathbf{s}^T \partial^2 \bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \gamma, \nu) \mathbf{s} > 0$$

537 for all \mathbf{s} satisfying (31) and (30). As $\limsup_{k \rightarrow \infty} w_k > 0$, by passing to a subsequence if necessary, we

538 assume $w_k \rightarrow \bar{w} > 0$. Then a similar analysis to the above proof gives

$$\mathbf{s}^T \partial^2 \bar{L}(\bar{\mathbf{x}}, \bar{\mathbf{y}}; \gamma, \nu) \mathbf{s} \geq \liminf w_k > 0.$$

539 □

540 C.4 Proof of Theorem 3.5

Proof. By Lemma D.9, we obtain that Assumption A.8 holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$. Then we can obtain that the sequence $\{(\mathbf{x}_{k',l}^*, \mathbf{y}_{k',l}^*)\}_{l=1}^\infty$, by passing to a subsequence if necessary, converges to $(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})$ by Lemma A.9. Furthermore, by Lemma A.10, by passing to a subsequence if necessary, we know

$$\|(\mathbf{x}_{k',l}^*, \mathbf{y}_{k',l}^*) - (\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})\| \leq \delta_{k'} \tau_{k',l},$$

541 where $\delta_{k'} = \sup_{l \geq 1} \frac{\|(\mathbf{x}_{k',l}^*, \mathbf{y}_{k',l}^*) - (\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})\|}{\tau_{k',l}}$ is bounded. We set $T_k \geq \frac{1}{2} + \log_\eta \frac{\delta_k}{\delta_{k-1}}$ in each iteration.

Let $\{w_{k'}\}$ be a subsequence of $\{w_k\}$ such that $\lim_{k' \rightarrow \infty} w_{k'} = \limsup_{k \rightarrow \infty} w_k$. As discussed right before this theorem, by Lemma A.9 and A.10 we have $\|(\mathbf{x}_{k',l}^*, \mathbf{y}_{k',l}^*) - (\bar{\mathbf{x}}_k', \bar{\mathbf{y}}_k')\| \leq \delta_{k',l} \tau_{k',l}$. Note that

$$\|(\mathbf{x}_{k',T_{k'}}^*, \mathbf{y}_{k',T_{k'}}^*) - (\bar{\mathbf{x}}, \bar{\mathbf{y}})\| \leq \|(\mathbf{x}_{k',T_{k'}}^*, \mathbf{y}_{k',T_{k'}}^*) - (\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})\| + \|(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}) - (\bar{\mathbf{x}}, \bar{\mathbf{y}})\|.$$

Note that

$$\|(\mathbf{x}_{k',T_{k'}}^*, \mathbf{y}_{k',T_{k'}}^*) - (\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})\| \leq \delta_{k'} \tau_{k',T_{k'}} = \delta_{k'} \frac{\tau_{k'-1,T_{k'-1}}}{\eta^{T_{k'}}} \leq \delta_{k'} \frac{\tau_{k'-1,T_{k'-1}}}{\sqrt{\eta} \frac{\delta_{k'}}{\delta_{k'-1}}} = \frac{\delta_{k'-1} \tau_{k'-1,T_{k'-1}}}{\sqrt{\eta}},$$

where the second inequality is due to (3) in the assumption. Using the above fact from k' to 2, we obtain

$$\|(\mathbf{x}_{k',T_{k'}}^*, \mathbf{y}_{k',T_{k'}}^*) - (\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})\| \leq \frac{\delta_1 \tau_{1,T_1}}{\eta^{\frac{k'-1}{2}}} \rightarrow 0 \quad (40)$$

as $k' \rightarrow \infty$. Because $\{(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})\} \rightarrow (\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we have $\|(\mathbf{x}_{k',T_{k'}}^*, \mathbf{y}_{k',T_{k'}}^*) - (\bar{\mathbf{x}}, \bar{\mathbf{y}})\| \rightarrow 0$.

Let $\gamma_k = (\lambda_1^k, \lambda_2^k, \dots, \lambda_{\kappa}^k)^T$, $\nu_k = \lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))$ as in the proof in Theorem 3.3. For notational simplicity, let $\mathbf{a} = (\mathbf{x}, \mathbf{y})$. As $F(\mathbf{a})$, $\mathbf{c}(\mathbf{x})$ and $\partial_{\mathbf{y}} f(\mathbf{a})$ are twice continuously differentiable, we have

$$\begin{aligned} & \partial \bar{L}(\bar{\mathbf{a}}_{k'}; \gamma_{k'}, \nu_{k'}) - \partial \bar{L}(\bar{\mathbf{a}}_{k'}; \gamma_{k'}, \nu_{k'}) \\ &= (\partial^2 F(\bar{\mathbf{a}}_{k'}) + \partial^2 (\mathbf{c}(\bar{\mathbf{x}}_{k'})^T \gamma_{k'}) + \partial^2 (\partial_{\mathbf{y}} f(\bar{\mathbf{a}}_{k'})^T \nu_{k'})) (\bar{\mathbf{a}}_{k'} - \mathbf{a}_{k',T_{k'}}^*) \end{aligned}$$

for some $\tilde{\mathbf{a}}_{k'}$ in the line $[\bar{\mathbf{a}}_{k'}, \mathbf{a}_{k',T_{k'}}^*]$, and there exists some absolute constants $C_1, C_2, C_3 > 0$ such that

$$\|\partial^2 F(\tilde{\mathbf{a}}_{k'})\| \leq C_1, \quad \|\partial^2 (\mathbf{c}(\bar{\mathbf{x}}_{k'})^T \gamma_{k'})\| \leq C_2 \|\gamma_{k'}\|, \quad \|\partial^2 (\partial_{\mathbf{y}} f(\tilde{\mathbf{a}}_{k'})^T \nu_{k'})\| \leq C_3 \|\nu_{k'}\|$$

for all $k' > 0$ as $\tilde{\mathbf{a}}_{k'}$ is bounded due to $\bar{\mathbf{a}}_{k'} \rightarrow \bar{\mathbf{a}}$ and $\mathbf{a}_{k',T_{k'}}^* \rightarrow \bar{\mathbf{a}}$. Hence due to $\gamma_{k'} \rightarrow \bar{\gamma}$ and $\nu_{k'} \rightarrow \bar{\nu}$, we have

$$\begin{aligned} & \|\partial \bar{L}(\bar{\mathbf{a}}_{k'}; \gamma_{k'}, \nu_{k'}) - \partial \bar{L}(\bar{\mathbf{a}}_{k'}; \bar{\gamma}, \bar{\nu})\| \\ & \leq O(C_1 + C_2 \|\bar{\gamma}\| + C_3 \|\bar{\nu}\|) \|\bar{\mathbf{a}}_{k'} - \mathbf{a}_{k',T_{k'}}^*\| \\ & = O(\|\bar{\mathbf{a}}_{k'} - \mathbf{a}_{k',T_{k'}}^*\|) \stackrel{(40)}{=} O(\eta^{-\frac{k'}{2}}). \end{aligned}$$

Therefore, combining the above fact with Theorem 3.3, we have

$$\|\partial \bar{L}(\mathbf{x}_{k',T_{k'}}^*, \mathbf{y}_{k',T_{k'}}^*; \gamma_{k'}, \nu_{k'})\| = O(\eta^{-\frac{k'}{2}}) + O(u^{-\frac{k'}{2}}).$$

555

□

D Affiliated Lemmas

Lemma D.1. Suppose $A \in \mathbb{R}^{n \times m}$ ($n \geq m$) is a matrix of full column rank, and there is a sequence of matrices A_k converging to A by element as $k \rightarrow \infty$. Then there exists $k_0 > 0$ such that for all $k > k_0$, A_k is of full column rank.

Proof. We utilize reduction to absurdity. Suppose that there exists a subsequence of A_k , denoted by $A_{k'}$, that $A_{k'}$ is not of full column rank. Then $\lim_{k' \rightarrow \infty} A_{k'} = A$. Since A is of full column rank, we know that the linear equation $A\mathbf{x} = \mathbf{0}$ has a unique solution $\mathbf{x} = \mathbf{0}$. On the other hand, the linear equations $A_{k'}\mathbf{x} = \mathbf{0}$ has non-zero solutions. We can choose a solution $\mathbf{x}_{k'}$ satisfying $\|\mathbf{x}_{k'}\| = 1$. So we can find a subsequence $\{\mathbf{x}_t\}$ satisfying $\lim_{t \rightarrow \infty} \mathbf{x}_t = \mathbf{x}_0$ as $\|\mathbf{x}_t\|$ is bounded, and absolutely $\|\mathbf{x}_0\| = \lim_{t \rightarrow \infty} \|\mathbf{x}_t\| = 1$. Then we have $A\mathbf{x}_0 = \lim_{t \rightarrow \infty} A_t \mathbf{x}_t = \mathbf{0}$, which contradicts $A\mathbf{x} = \mathbf{0}$ has unique solution $\mathbf{x} = \mathbf{0}$. □

Lemma D.2. Suppose A_k and A are defined as in Lemma D.1, $\{\mathbf{b}_k \in \mathbb{R}^m\}_{k=1}^{\infty}$ is a vector sequence and $\{\|A_k \mathbf{b}_k\|\}_{k=1}^{\infty}$ is bounded, then $\{\|\mathbf{b}_k\|\}$ is bounded.

Proof. As $\{\|A_k \mathbf{b}_k\|\}_{k=1}^{\infty}$ is bounded, there exists some $M_0 > 0$ such that $\|A_k \mathbf{b}_k\| = \sqrt{\mathbf{b}_k^T A_k^T A_k \mathbf{b}_k} < M_0$ for all $k > 0$. By Lemma D.1 we have known that there exists $k_0 > 0$ such that for all $k > k_0$, A_k is of full column matrix. Therefore, $A_k^T A_k$ is positive definite when

572 k is large enough. On the other hand, $A^T A \succ 0$ because A is of full column rank. We define the
 573 minimum eigenvalue of $A^T A$ as σ_1 , which is positive. As the minimum eigenvalue of a matrix is
 574 continuous w.r.t. its elements, we obtain that there exists a positive integer K such that for all
 575 $k > K$, $A_k^T A_k \succeq \frac{\sigma_1}{2} I$. As $\|A_k \mathbf{b}_k\| < M_0$, it follows that $\|\mathbf{b}_k\| \leq \sqrt{\frac{2}{\sigma_1}} M_0$ when k is large enough.
 576 Thus $\|\mathbf{b}_k\|$ is bounded when k is large enough. \square

577 **Lemma D.3.** Suppose A_k and A are defined as in Lemma D.1. Then for every vector \mathbf{p} satisfying
 578 $A^T \mathbf{p} = 0$, there exists a sequence of vectors denoted as $\{\mathbf{p}_k\}_{k=1}^\infty$ converging to \mathbf{p} and $A_k^T \mathbf{p}_k = 0$.

579 *Proof.* Note that $\lim_{k \rightarrow \infty} A_k^T \mathbf{p} = A^T \mathbf{p} = 0$. Let \mathbf{p}_k be such that $A_k^T (\mathbf{p}_k - \mathbf{p}) = -A_k^T \mathbf{p}$, which is
 580 equivalent to $A_k^T \mathbf{p}_k = 0$. Since when k is large enough $A_k^T A_k$ is invertible, we can let $\mathbf{p}_k - \mathbf{p} =$
 581 $-A_k (A_k^T A_k)^{-1} A_k^T \mathbf{p}$. In the proof of Lemma D.2 we have shown there exists a $\sigma_1 > 0$ such
 582 that $A_k^T A_k \geq \frac{\sigma_1}{2} I$ if k is large enough, then $(A_k^T A_k)^{-1} \leq \frac{2}{\sigma_1} I$ when k is sufficiently large. So
 583 $\|A_k (A_k^T A_k)^{-1} A_k^T \mathbf{p}\| = \sqrt{\mathbf{p}^T A_k (A_k^T A_k)^{-1} A_k^T \mathbf{p}} \leq \sqrt{\frac{2}{\sigma_1}} \|A_k^T \mathbf{p}\|$, which implies $\lim_{k \rightarrow \infty} \mathbf{p}_k - \mathbf{p} =$
 584 $\lim_{k \rightarrow \infty} -A_k (A_k^T A_k)^{-1} A_k^T \mathbf{p} = 0$. By selecting $\mathbf{p}_k = \mathbf{p} - A_k (A_k^T A_k)^{-1} A_k^T \mathbf{p}$ when k is large enough,
 585 we finish the proof. \square

586 **Lemma D.4.** Consider the LL function $f(\mathbf{x}, \mathbf{y}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$. When $\mathbf{z}^*(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is
 587 a continuous function of \mathbf{x} in a neighborhood of $\bar{\mathbf{x}}$, then we can find a neighborhood of $\bar{\mathbf{x}}$ in which
 588 $\mathbf{z}^*(\mathbf{x})$ is a continuously differentiable function of \mathbf{x} .

589 *Proof.* We have known that $\bar{\mathbf{y}} \in \operatorname{argmin}_{\mathbf{y}} f(\bar{\mathbf{x}}, \mathbf{y})$ from Lemma C.1, and thus $\partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = 0$. Since
 590 we have assumed $\partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \succ 0$, thus $\partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ is invertible. By the implicit function theorem,
 591 there exist open sets $U \subseteq \mathbb{R}^{m+n}$ and $W \subseteq \mathbb{R}^m$ with $(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in U$ and $\bar{\mathbf{x}} \subseteq W$, satisfying that for
 592 every $\mathbf{x} \in W$, there exists a unique \mathbf{y} such that

$$(\mathbf{x}, \mathbf{y}) \in U \quad \text{and} \quad \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = 0$$

593 and \mathbf{y} can be seen as a continuously differentiable function of \mathbf{x} in W . We denote it as $\mathbf{y} = \mathbf{h}(\mathbf{x})$, $\mathbf{x} \in$
 594 W . Next we need to prove $\mathbf{z}^*(\mathbf{x}) = \mathbf{h}(\mathbf{x})$ near $\bar{\mathbf{x}}$.
 595 Since $\mathbf{z}^*(\mathbf{x})$ is continuous w.r.t. \mathbf{x} in a neighborhood of $\bar{\mathbf{x}}$ and $\mathbf{z}^*(\bar{\mathbf{x}}) = \bar{\mathbf{y}}$, we can find open sets
 596 $V_1 \subseteq W$ and $V_2 \subseteq \mathbb{R}^n$ such that $\bar{\mathbf{x}} \in V_1$, $\bar{\mathbf{y}} \in V_2$, $V_1 \times V_2 \subseteq U$, and

$$\mathbf{z}^*(\mathbf{x}) \in V_2 \quad \text{for all } \mathbf{x} \in V_1.$$

597 However, $\mathbf{z}^*(\mathbf{x})$ satisfies $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) = 0$. Then by the uniqueness of \mathbf{y} when $(\mathbf{x}, \mathbf{y}) \in U$, we
 598 have $\mathbf{z}^*(\mathbf{x}) = \mathbf{h}(\mathbf{x})$ in V_1 . Noting that \mathbf{h} is continuously differentiable, we complete the proof. \square

599 **Lemma D.5.** Under Assumption 3.1, there exists $k_0 > 0$ such that for all $k > k_0$, we have $\mathcal{A}_k \subseteq \bar{\mathcal{A}}$
 600 and LICQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for problem (6).

601 *Proof.* Recall $\mathcal{A}_k = \{i : c_i(\bar{\mathbf{x}}_k) = 0, i = 1, 2, \dots, s\}$ and $\bar{\mathcal{A}} = \{i : c_i(\bar{\mathbf{x}}) = 0, i = 1, 2, \dots, s\}$. We
 602 claim that there exists $k_0 > 0$ such that when $k > k_0$, $\mathcal{A}_k \subseteq \bar{\mathcal{A}}$. If not, we can find an index $i_0 \notin \bar{\mathcal{A}}$
 603 but $c_{i_0}(\bar{\mathbf{x}}_k) = 0$ for a subsequence of $\{\bar{\mathbf{x}}_k\}$. Then we have $c_{i_0}(\bar{\mathbf{x}}) = 0$, and thus $i_0 \in \bar{\mathcal{A}}$, which
 604 contradicts the assumption.
 605 Let $\bar{\mathcal{A}} = \{i_1, i_2, \dots, i_t\}$. Then we have

$$J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) = (\partial c_{i_1}(\bar{\mathbf{x}}), \partial c_{i_2}(\bar{\mathbf{x}}), \dots, \partial c_{i_t}(\bar{\mathbf{x}}))^T \quad \text{and} \quad J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k) = (\partial c_{i_1}(\bar{\mathbf{x}}_k), \partial c_{i_2}(\bar{\mathbf{x}}_k), \dots, \partial c_{i_t}(\bar{\mathbf{x}}_k))^T.$$

606 Since $\mathcal{A}_k \subseteq \bar{\mathcal{A}}$ when k is large enough, $J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k)$ is a submatrix of $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k)$. Note that $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}})^T$ is of
 607 full column rank by Assumption 3.1. By Lemma D.1, we know there exists $k_0 > 0$ such that for all
 608 $k > k_0$, $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k)^T$ is of full column rank. Then $J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k)^T$ is of full column rank because $\mathcal{A}_k \subseteq \bar{\mathcal{A}}$.
 609 When c_0 is inactive at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$, i.e. $c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \mu_k \neq 0$, LICQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for (6) because
 610 $J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k)^T$ is of full column rank.

611 Now let us consider the case that c_0 is active, i.e. $c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = \mu_k$. We claim that we cannot have
 612 $\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = \mathbf{0}$ when k is sufficiently large. Since $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \succ 0$, $\mathbf{z}^*(\mathbf{x})$ is continuous in a
 613 neighborhood of $\bar{\mathbf{x}}$ and $\lim_{k \rightarrow \infty} (\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = (\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we obtain that when k is large enough, $\bar{\mathbf{y}}_k$ is in
 614 a neighborhood of $\mathbf{z}^*(\bar{\mathbf{x}}_k)$ where f is strongly convex w.r.t. \mathbf{y} . Then $\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = \mathbf{0}$ implies

615 $\bar{\mathbf{y}}_k = \mathbf{z}^*(\bar{\mathbf{x}}_k)$. This yields $c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = 0 < \mu_k$, which contradicts the fact that $c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = \mu_k$ is
 616 inactive. So we must have $\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \neq \mathbf{0}$. This, together with the fact that $J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k)^T$ is of full
 617 column rank and $\partial_{\mathbf{y}} c_i(\mathbf{x}) = \mathbf{0} \forall i \in [\kappa]$, yields LICQ for (6). \square

618 **Lemma D.6.** *Under Assumption 3.1, LICQ holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ for (10).*

619 *Proof.* It suffices to show the transpose of the Jacobian

$$J_0 \triangleq \begin{pmatrix} (\partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^T & \partial_{c_{i_1}}(\bar{\mathbf{x}}) & \cdots & \partial_{c_{i_t}}(\bar{\mathbf{x}}) \end{pmatrix} = \begin{pmatrix} (\partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^T & \partial_{\mathbf{x}} c_{i_1}(\bar{\mathbf{x}}) & \cdots & \partial_{\mathbf{x}} c_{i_t}(\bar{\mathbf{x}}) \\ (\partial_{\mathbf{y}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^T & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix}$$

620 is of full column rank, where $\bar{\mathcal{A}} = \{i_1, i_2, \dots, i_t\}$ is the index set for constraints $c_i(\mathbf{x}) \forall i \in [\kappa]$. This
 621 is true as

$$\text{rank}(J_0) \geq \text{rank}(\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}})) + \text{rank}([\partial_{c_{i_1}}(\bar{\mathbf{x}}, \bar{\mathbf{y}}), \dots, \partial_{c_{i_t}}(\bar{\mathbf{x}}, \bar{\mathbf{y}})]) = n + t.$$

622 \square

623 **Lemma D.7.** *Given the same assumptions and notation in the proof of Theorem 3.2, $\{\lambda_i^k\}_{k=1}^\infty \forall i \in [\kappa]$
 624 are bounded.*

625 *Proof.* In the proof of Theorem 3.2, we show that $\|\sum_{i=1}^\kappa \lambda_i^k \partial_{\mathbf{x}} c_i(\bar{\mathbf{x}}_k)\|$ is uniformly bounded for k
 626 in (22). We also show $\mathcal{A}_k \subseteq \bar{\mathcal{A}}$ when k is large enough in the proof of Lemma D.5. Since $\lambda_i^k = 0$
 627 if $i \notin \mathcal{A}_k$ due to complementary slackness in the KKT conditions, we know $\|J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k)^T \boldsymbol{\lambda}_{\bar{\mathcal{A}}}^k\|$ is
 628 uniformly bounded where $\boldsymbol{\lambda}_{\bar{\mathcal{A}}}^k$ is an acceptable Lagrange multiplier of the index set $\bar{\mathcal{A}}$. Note that
 629 $\lim_{k \rightarrow \infty} J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k)^T = J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}})^T$ and that $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}})^T$ is of full column rank. By Lemma D.2, we deduce the
 630 boundedness of $\{\lambda_i^k\}$ for $i \in \bar{\mathcal{A}}$. This, together with $\lambda_i^k = 0$ if $i \notin \bar{\mathcal{A}}$, completes the proof. \square

631 The following Lemmas are for Theorem 3.4 and 3.5, some proof of them are base on the proof of
 632 Theorem 3.2.

633 **Lemma D.8.** *Under Assumption 3.1, if additionally strict complementarity holds at $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ and
 634 $\partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \neq \mathbf{0}$, then strict complementarity holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ when k is large enough.*

635 *Proof.* Note that when k is large enough, $\mathcal{A}_k \subseteq \bar{\mathcal{A}}$ and LICQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for μ_k relaxation
 636 problem by Lemma D.5, then the acceptable Lagrangian multiplier would be unique at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$. If
 637 there exists a subsequence $\{\lambda_i^k\}$ such that $\lambda_i^k = 0$, then in (23), λ_i would be 0, which contradicts the
 638 assumption. Therefore, when k is large enough, $\lambda_i^k > 0$ for $i \in \mathcal{A}_k$.
 639 Now we show $\lambda_0^k \neq 0$ when k is sufficiently large, which would finish the proof. Since we have
 640 assumed $\partial_{\mathbf{y}} F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \neq \mathbf{0}$, then $\partial_{\mathbf{y}} F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \neq \mathbf{0}$ when k is large enough. By (15), $\lambda_0^k \neq 0$ holds. \square

641 **Lemma D.9.** *Suppose that the condition of Theorem 3.4 holds with $\limsup_{k \rightarrow \infty} w_k > 0$, and
 642 additionally strict complementarity holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ when k is large enough. Then Assumption A.8
 643 holds at $(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})$ when k' is sufficiently large.*

644 *Proof.* Since MFCQ is weaker than LICQ and LICQ holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for problem (6) by Lemma
 645 D.5 when k is large enough, we only need to show that (2) holds in Assumption A.8. To begin with,
 646 because of LICQ, the acceptable Lagrangian multiplier would be unique at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ when k is large
 647 enough.

648 We firstly consider the situation that $c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}) \leq \mu_{k'}$ is always active when k' is large enough. Let
 649 $\boldsymbol{\lambda}_{\mathcal{A}_{k'}}^{k'}$ be the acceptable Lagrange multiplier of index set $\mathcal{A}_{k'}$. Note that the KKT conditions of (6)
 650 implies

$$\nabla L(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}, \lambda_0^{k'}, \boldsymbol{\lambda}^{k'}) = \partial F(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}) + \lambda_0^{k'} \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}) + \sum_{i=1}^\kappa \lambda_i^{k'} \partial c_i(\bar{\mathbf{x}}_{k'}) = \mathbf{0}. \quad (41)$$

651 For any \mathbf{p} satisfying $\partial F(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \mathbf{p} = 0$ and $\begin{pmatrix} J_{\mathcal{A}_{k'}}(\bar{\mathbf{x}}_{k'}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \end{pmatrix} \mathbf{p} \geq \mathbf{0}$, we have

652 $\begin{pmatrix} \boldsymbol{\lambda}_{\mathcal{A}_{k'}}^{k'} \\ \lambda_0^{k'} \end{pmatrix}^T \begin{pmatrix} J_{\mathcal{A}_{k'}}(\bar{\mathbf{x}}_{k'}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \end{pmatrix} \mathbf{p} = 0$ by (41). Then, using $\boldsymbol{\lambda}_{\mathcal{A}_{k'}}^{k'} > \mathbf{0}$ and $\lambda_0^{k'} > 0$ due to

strict complementarity, we have $\begin{pmatrix} J_{\mathcal{A}_{k'}}(\bar{\mathbf{x}}_{k'}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \end{pmatrix} \mathbf{p} = 0$. Moreover, $\partial F(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \mathbf{p} = 0$ and $\begin{pmatrix} J_{\mathcal{A}_{k'}}(\bar{\mathbf{x}}_{k'}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \end{pmatrix} \mathbf{p} \geq 0$ are equivalent to $\begin{pmatrix} J_{\mathcal{A}_{k'}}(\bar{\mathbf{x}}_{k'}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \end{pmatrix} \mathbf{p} = \mathbf{0}$ from (41). Now it suffices to show for all \mathbf{p} satisfying $\begin{pmatrix} J_{\mathcal{A}_{k'}}(\bar{\mathbf{x}}_{k'}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'})^T \end{pmatrix} \mathbf{p} = \mathbf{0}$, we must have $\mathbf{p}^T \partial^2 L(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}, \lambda_{k'}) \mathbf{p} > 0$. Indeed, since when k' is large enough we have $w_{k'} > 0$, by second-order necessary conditions at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ and the definition of w_k , we have $\mathbf{p}^T \partial^2 L(\bar{\mathbf{x}}_{k'}, \bar{\mathbf{y}}_{k'}, \lambda_{k'}) \mathbf{p} \geq w_{k'} \|\mathbf{p}\|^2 > 0$. On the other hand, when there is a subsequence of $\{k'\}$, denoted by $\{k'_j\}$, such that constraint $c_0(\bar{\mathbf{x}}_{k'_j}, \bar{\mathbf{y}}_{k'_j}) < \mu_{k'_j}$. We can prove the results in a similar way. In this situation $\lambda_0^{k'_j} = 0$, and the main process is to prove $\partial F(\bar{\mathbf{x}}_{k'_j}, \bar{\mathbf{y}}_{k'_j})^T \mathbf{p} = 0$ and $\begin{pmatrix} J_{\mathcal{A}_{k'_j}}(\bar{\mathbf{x}}_{k'_j}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'_j}, \bar{\mathbf{y}}_{k'_j})^T \end{pmatrix} \mathbf{p} \geq 0$ are equivalent to $\begin{pmatrix} J_{\mathcal{A}_{k'_j}}(\bar{\mathbf{x}}_{k'_j}) & 0 \\ \partial c_0(\bar{\mathbf{x}}_{k'_j}, \bar{\mathbf{y}}_{k'_j})^T \end{pmatrix} \mathbf{p} = 0$ under strict complementarity. We omit the proof for simplicity. \square

In the following lemma, we show that the SONCs holds for problem (6) at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ if k is sufficiently large.

Lemma D.10. *Given the same assumptions and notation in the proof of Theorem 3.4, we have*

$$\mathbf{s}^T \partial^2 F(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{s} - \sum_{i=1}^{\kappa} \lambda_i \mathbf{s}^T \partial^2 c_i(\bar{\mathbf{x}}) \mathbf{s} + \sum_{i=1}^n q_i \mathbf{x}^T \partial_{y_i} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{l}_x \geq \liminf_{k \rightarrow \infty} w_k \|\mathbf{s}\|^2, \quad (42)$$

where $\boldsymbol{\lambda} = \lim_{k \rightarrow \infty} \boldsymbol{\lambda}_i^k$ and $\mathbf{q} = \lim_{k \rightarrow \infty} \lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))$.

Proof. Recall that

$$P(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \begin{pmatrix} I_m & 0_{m \times n} \\ (\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}, \bar{\mathbf{y}}))^{-1} \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) & I_n \end{pmatrix}.$$

The outline of the proof is that we will show for each pair (\mathbf{s}, \mathbf{l}) satisfying $\mathbf{s} = P(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{l}$, $J_{\mathcal{A}}(\bar{\mathbf{x}}) \mathbf{l}_x = 0$ and $\mathbf{l}_y = 0$, there exists convergent sequences $\mathbf{s}_k \in \mathbb{R}^{m+n}$ and $\mathbf{l}_k \in \mathbb{R}^{m+n}$ satisfying $\mathbf{s}_k = P(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{l}_k$, $\mathbf{s}_k \rightarrow \mathbf{s}$ and $\mathbf{l}_k \rightarrow \mathbf{l} \triangleq (\mathbf{l}_x^T, \mathbf{l}_y^T)^T$, such that we can find a sequence $\{\mathbf{s}_k\}_{k=1}^{\infty}$ that satisfies $\partial c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^T \mathbf{s}_k = 0$ and $\partial c_i(\bar{\mathbf{x}}_k)^T \mathbf{s}_k = 0$ for $i \in [\kappa]$, and (42) holds. For problem (6), the Lagrange function, $\partial^2 L_k(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \lambda_0^k, \boldsymbol{\lambda}^k)$ has the form,

$$\partial^2 L_k(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \lambda_0^k, \boldsymbol{\lambda}^k) = \partial^2 F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \sum_{i=1}^{\kappa} \lambda_i^k \partial^2 c_i(\bar{\mathbf{x}}_k) + \lambda_0^k \partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k).$$

In the following, we will show the SONCs holds for problem (6) by finding a sequence \mathbf{s}_k satisfying Definition A.4. By Proposition 2.2, we have

$$\partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = \begin{pmatrix} \partial_{\mathbf{x}\mathbf{x}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}\mathbf{x}}^2 f^*(\bar{\mathbf{x}}_k) & \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \\ \partial_{\mathbf{x}} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) & \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \end{pmatrix}.$$

Let $P_k = P(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$. Then one may check the following holds

$$\partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = P_k^T \begin{pmatrix} Q(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - Q(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) & 0_{m \times n} \\ 0_{n \times m} & \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \end{pmatrix} P_k, \quad (43)$$

using the expression of $\partial_{\mathbf{x}\mathbf{x}}^2 f_{\mu}^*(\mathbf{x})$ in Table 1. By the third order continuously differentiability of f , $(\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1}$ and thus $Q(\mathbf{x}, \mathbf{y})$ are a continuously differentiable function due to Theorem 8.3 in Magnus & Neudecker [16]. Now we apply Taylor expansion to every entry y_i of the matrix $Q(\mathbf{x}, \mathbf{y})$.

$$Q(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = Q(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) + \sum_{i=1}^n (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))_i \partial_{y_i} Q(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) + W_k, \quad (44)$$

where W_k satisfies $\|W_k\| = o(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|)$, and y_k^i is the i th entry of $\bar{\mathbf{y}}_k$. Let $\mathbf{l}_k = P_k \mathbf{s}_k$ and $\mathbf{l}_k = ((\mathbf{l}_x^k)^T, (\mathbf{l}_y^k)^T)^T$. Then using (43) and (44), we have

$$\mathbf{s}_k^T \partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{s}_k = \sum_{i=1}^n (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))_i (\mathbf{l}_x^k)^T \partial_{y_i} Q(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) \mathbf{l}_x^k + (\mathbf{l}_y^k)^T \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{l}_y^k + o(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| \|\mathbf{l}_x^k\|^2). \quad (45)$$

Let \mathbf{s}_k be such that

$$\partial c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^T \mathbf{s}_k = 0 \text{ and } \partial c_i(\bar{\mathbf{x}}_k)^T \mathbf{s}_k = 0 \forall i \in [\kappa]. \quad (46)$$

Then we have $\partial c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^T P_k^{-1} \mathbf{l}_k = 0$, and thus

$$\left(\partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k) - \frac{\partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^{-1} \partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)}{\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)} \right)^T \begin{pmatrix} \mathbf{l}_x^k \\ \mathbf{l}_y^k \end{pmatrix} = 0, \quad (47)$$

due to the definition of P_k and the expression of ∂c_0 . Using Taylor expansion, we have

$$\partial_{\mathbf{x}} f^*(\bar{\mathbf{x}}_k) = \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) = \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) + \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) (\mathbf{z}^*(\bar{\mathbf{x}}_k) - \bar{\mathbf{y}}_k) + \mathbf{v}_1^k$$

where \mathbf{v}_1^k satisfies $\|\mathbf{v}_1^k\| = O(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|^2)$, and, due to $\partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) = \mathbf{0}$ by definition,

$$\mathbf{0} = \partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) = \partial_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) + \partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) (\mathbf{z}^*(\bar{\mathbf{x}}_k) - \bar{\mathbf{y}}_k) + \mathbf{v}_2^k,$$

where \mathbf{v}_2^k satisfies $\|\mathbf{v}_2^k\| = O(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|^2)$. By substituting the above two equations into (47), we obtain

$$-(\mathbf{l}_x^k)^T \mathbf{v}_1^k + (\mathbf{l}_x^k)^T \partial_{\mathbf{y}} \partial_{\mathbf{x}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^{-1} \mathbf{v}_2^k - (\mathbf{l}_y^k)^T \partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) (\mathbf{z}^*(\bar{\mathbf{x}}_k) - \bar{\mathbf{y}}_k) + (\mathbf{l}_y^k)^T \mathbf{v}_2^k = 0,$$

which further implies

$$|(\mathbf{l}_y^k)^T \partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))| = O(\|\mathbf{l}_x^k\| \cdot \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|^2).$$

When k is large enough, we have $\partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) \succ aI$ for some $a > 0$ due to (1) and (2) in Assumption 3.1.

Note that $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k)$ converges to $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}})$ as $\bar{\mathbf{x}}_k \rightarrow \bar{\mathbf{x}}$. By using Lemma D.3, we know for every \mathbf{l}_x satisfying $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}) \mathbf{l}_x = 0$, we can find $\{\mathbf{l}_x^k\}$ satisfying $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k) \mathbf{l}_x^k = \mathbf{0}$ and $\mathbf{l}_x^k \rightarrow \mathbf{l}_x$. Now we select \mathbf{l}_x^k such that $\mathbf{l}_x^k \rightarrow \mathbf{l}_x$. Let \mathbf{l}_y^k be parallel to $\partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))$. We must have

$$\|\mathbf{l}_y^k\| = O(\|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\|). \quad (48)$$

By (45) we have

$$\begin{aligned} \lambda_0^k \mathbf{s}_k^T \partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{s}_k &= \lambda_0^k \sum_{i=1}^n (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))_i (\mathbf{l}_x^k)^T \partial_{y_i} Q(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k)) \mathbf{l}_x^k + \lambda_0^k (\mathbf{l}_y^k)^T \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{l}_y^k \\ &\quad + o(\lambda_0^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| \|\mathbf{l}_x^k\|^2). \end{aligned} \quad (49)$$

In the proof of Theorem 3.2, we have shown that

$$\lambda_i^k \|\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)\| = O(1). \quad (50)$$

By passing to a subsequence if necessary, we can further assume $\lim_{k \rightarrow \infty} \lambda_0^k (\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k)) \rightarrow \bar{\mathbf{q}}$. From (49) we have

$$\lim_{k \rightarrow \infty} \lambda_0^k \mathbf{s}_k^T \partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{s}_k = \sum_{i=1}^n q_i \mathbf{l}_x^T \partial_{y_i} Q(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mathbf{l}_x, \quad (51)$$

where the limit follows from (48), (50), $\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{y}}_k) \rightarrow \mathbf{0}$ and

$$\|\lambda_0^k (\mathbf{l}_y^k)^T \partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{l}_y^k\| \leq \|\lambda_0^k \mathbf{l}_y^k\| \cdot \|\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)\| \cdot \|\mathbf{l}_y^k\| \rightarrow 0.$$

Recall that we have chosen $J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k) \mathbf{l}_x^k = \mathbf{0}$. For any sufficiently large k , we have $J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k) \mathbf{l}_x^k = \mathbf{0}$ as we have shown $\bar{\mathcal{A}} \supseteq \mathcal{A}_k$ in Lemma D.5. Let \mathbf{l}_y^k be such that (47) holds. Due to (48) and $\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{y}}_k) \rightarrow \mathbf{0}$, we have $\lim_{k \rightarrow \infty} \mathbf{l}_y^k = \mathbf{0}$. Correspondingly, $\lim_{k \rightarrow \infty} \mathbf{l}_k = \mathbf{l}$. Using $\lim_{k \rightarrow \infty} P(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = P(\bar{\mathbf{x}}, \bar{\mathbf{y}})$, we then have $\lim_{k \rightarrow \infty} \mathbf{s}_k = \mathbf{s}$. We remark the well definedness of \mathbf{l}^k and thus \mathbf{s}_k . Note that by the definition of $P(\mathbf{x}, \mathbf{y})$, $\mathbf{s}_k = P(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{l}_k$ and (46), we see that $\partial c_i(\bar{\mathbf{x}}_k)^T \mathbf{s}_k = 0 \forall i \in [\kappa]$ is equivalent to

$$(J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k) \mathbf{0}_{t \times n}) \mathbf{s}_k = (J_{\mathcal{A}_k}(\bar{\mathbf{x}}_k) \mathbf{0}_{t \times n}) \mathbf{l}_k = J_{\bar{\mathcal{A}}}(\bar{\mathbf{x}}_k) \mathbf{l}_x^k = \mathbf{0}.$$

And $\lim_{k \rightarrow \infty} \mathbf{I}_k = \mathbf{I}$ can be implied by (47), which is just $\partial c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)^T \mathbf{s}_k = 0$ in (46). The fact \mathbf{I}_k^k be parallel to $\partial_{\mathbf{y}\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{z}^*(\bar{\mathbf{x}}_k))(\bar{\mathbf{y}}_k - \mathbf{z}^*(\bar{\mathbf{x}}_k))$ follows from that \mathbf{I}_k^k is only involved in and does not violate the equation (47). That is, the well definedness of \mathbf{I}^k follows from the linear system (46).

Since the LICQ for (6) holds due to Lemma D.5, we obtain that the SONCs of (6) holds, thanks to Theorem A.6. We point out that as strict complementarity holds at $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ for sufficiently large k , we always have $\lambda_{\mathcal{A}_k} > 0$. Then by definition, we have the critical cone

$$\mathcal{C}((\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k), (\lambda_0^k, \boldsymbol{\lambda}^k)) \supseteq \{\mathbf{s} : \partial c_0(\bar{\mathbf{x}}_k)^T \mathbf{s} = 0, \partial c_i(\bar{\mathbf{x}}_k)^T \mathbf{s} = 0 \text{ for all } i \in \bar{\mathcal{A}}\}.$$

Hence for the \mathbf{s}_k satisfying (46) we always have $w_k \geq 0$ such that

$$\mathbf{s}_k^T \partial^2 L_k(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k; \lambda_0^k, \boldsymbol{\lambda}^k) \mathbf{s}_k \geq w_k \|\mathbf{s}_k\|^2,$$

which is equivalent to

$$\mathbf{s}_k^T \partial^2 F(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{s}_k - \sum_{i=1}^K \lambda_i^k \mathbf{s}_k^T \partial^2 c_i(\bar{\mathbf{x}}_k) \mathbf{s}_k + \lambda_0^k \mathbf{s}_k^T \partial^2 c_0(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \mathbf{s}_k \geq w_k \|\mathbf{s}_k\|^2.$$

Substituting (51) to the above inequality and taking the limit, we have the desired result. \square

E Supplementary Experiments

E.1 Bi-level Logistic Regression

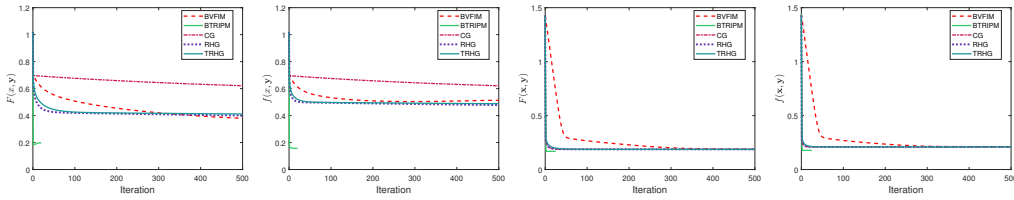


Figure 3: Comparison of BVFIM, RHG, TRHG, CG with BTRIPM for solving bilevel logistic regression based on iterations. The left two are on WIL data set and the right two are on CM data set.

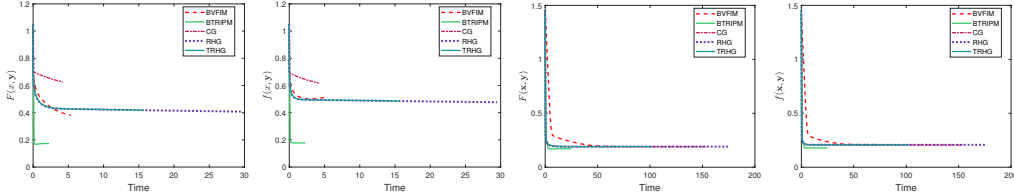


Figure 4: Comparison of BVFIM, RHG, TRHG, CG with BTRIPM for solving bilevel logistic regression based on time. The left two are on WIL data set and the right two are on CM data set.

In order to affirm the universality of our algorithm in various BLO problems, we attempt to address more complex issues in practical background with different functions. Similar as in Pedregosa [19], we test our algorithm on classification tasks with logistic regression models on two real-world datasets. For completeness, we add another gradient-based method for comparison as in Shaban et al. [22], which solves the lower-level problem through iterative optimization procedure with hyper gradient computed by truncated back-propagation. We denote it as TRHG and apply to optimize hyperparameters controlling the lower-level optimization procedure.

We partition each initial dataset into two sets: a train set $S_{\text{train}} = \{(\mathbf{b}_i, a_i)\}_{i=1}^{D_{\text{tr}}}$ and a test set $S_{\text{test}} = \{(\mathbf{d}_i, c_i)\}_{i=1}^{D_{\text{te}}}$. Here \mathbf{b}_i and \mathbf{d}_i denote the input features, and a_i and c_i denote the labels. To classify these labels, we need to estimate a regularization parameter in the widely used l_2 -regularized logistic regression model while the validation loss is defined as the logistic loss. Specifically, the problem is written as

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}, \mathbf{y} \in \mathbb{R}^n} \quad & \sum_{i \in S_{\text{test}}} \psi(c_i \mathbf{d}_i^T \mathbf{y}) \\ \text{s.t.} \quad & \mathbf{y} \in \arg\min_{\mathbf{y}} \sum_{i \in S_{\text{train}}} \psi(a_i \mathbf{b}_i^T \mathbf{y}) + e^x \|\mathbf{y}\|^2, \end{aligned}$$

Table 4: Comparison of BVFIM, RHG, TRHG, CG with BTRIPM in the aspect of the total time(s)

Dataset	BVFIM	CG	RHG	TRHG	BTRIPM
WIL	5.1769	4.2529	29.5356	15.7332	2.5402
CM	77.1380	164.2327	180.4677	107.8747	26.0561

where ψ is the logistic loss, i.e., $\psi(t) = \ln(1 + e^{-t})$. The first data set "Wireless Indoor Localization" (WIL)⁴ was collected to perform experimentation on how wifi signal strengths can be used to determine one of the indoor locations in [21]. 7 attributes in 2,000 instances are in data set and each is wifi signal strength observed on smartphone, which results in missing values as response variable. Another data set "Crowdsourced Mapping" (CM)⁵ was derived from two geospatial data sources, landsat time-series satellite imagery and crowdsourced georeferenced polygons with land cover labels [10]. In traditional cognitive, logistic regression is applied in binary classification. In our experiments, labels to classify are the land cover class resulted from 28 attributes in 10546 instances, whose terminal values are split into two categories to be transformed into a binary classification problem.

Figures E.1 reports classification results in the datasets WIL and CM, respectively. Both figures show BTRIPM achieves high accuracy in less iterations than other methods. BTRIPM appears more precise in the two datasets, while all the other algorithms find bad solutions for dataset WIL. In both datasets, BTRIPM only requires less than ten outer iterations to converge, while others need more than five hundred iterations. Table 4 further shows that the total time consumed by different algorithms on two different data sets. In general, BTRIPM expends less time than all the first-order algorithms.

E.2 Bi-level Multinomial Logistic Regression

Moreover, it is urgent to manifest the efficiency of our algorithm on large scale problems and solve instances of the bilevel problem with variables in a higher dimension. We now compare algorithms on multinomial logistic regression on a text application dataset.

The dataset "20 Newsgroup" was collected with approximately 20,000 newsgroup documents and organized into 20 different newsgroups, each corresponding to a different topic. For feasibility and convenience, the updated dataset sorted by date with duplicates and some headers removed has been provided and available⁶, which contains $N = 18,846$ text documents. We employ a processed version of updated dataset easy to read into Matlab. The features of text documents are rearranged into sparse matrix, in which each row represents an instance. Meanwhile, the labels are stored as a column vector with the same indexes.

We divide the data randomly into three equal segments, each with $N/3$ documents: a train set $\{X_{tr}, \mathbf{y}_{tr}\}$, a validation set $\{X_{val}, \mathbf{y}_{val}\}$ and a test set $\{X_{te}, \mathbf{y}_{te}\}$. Here X_{tr}, X_{val}, X_{te} are sparse matrix of features with sparsity as approximately 0.5%, and $\mathbf{y}_{tr}, \mathbf{y}_{val}, \mathbf{y}_{te}$ denote the labels. We aim to execute classification for labels with l_2 -regularized multinomial logistic regression model and define the validation loss as the cross-entropy loss in the following bilevel problem

$$\begin{aligned} \min_{\lambda \in \mathbb{R}^p} \quad & \text{CE}(X_{val}w(\lambda), \mathbf{y}_{val}) \\ \text{s.t.} \quad & w(\lambda) \in \underset{w \in \mathbb{R}^{p \times c}}{\text{argmin}} \text{CE}(X_{tr}w(\lambda), \mathbf{y}_{tr}) + \frac{1}{2cp} \sum_{i=1}^c \sum_{j=1}^p \exp(\lambda_j) w_{ij}^2 \end{aligned} \quad (52)$$

where CE is the average cross-entropy loss, $c = 20$ and $p = 26,214$. In our experiments, we estimate a regularization parameter and compute corresponding accuracy in the test set.

Table 5 and Figure E.2 report classification results in the dataset "20 Newsgroup". Among the operated algorithms, BTRIPM appears more precise and with less time consumption.

⁴<https://archive.ics.uci.edu/ml/datasets/Wireless+Indoor+Localization>

⁵<https://archive.ics.uci.edu/ml/datasets/Crowdsourced+Mapping>

⁶<http://qwone.com/~jason/20Newsgroups/>

Table 5: Comparison of BVFIM, RHG, TRHG, CG with BTRIPM in the aspect of accuracy and the total time (s).

Algorithm	Accuracy	Total time (s)
BVFIM	0.7690	1601.5
CG	0.7934	3159.5
RHG	0.7934	2393.4
TRHG	0.7929	1618.1
BTRIPM	0.8155	1022.6

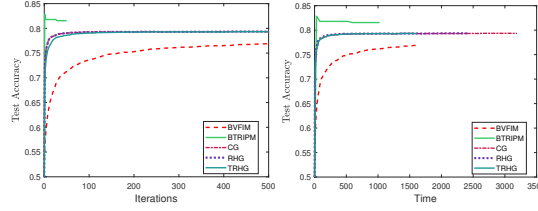


Figure 5: Comparison of BVFIM, RHG, TRHG, CG with BTRIPM for solving Bi-level Multinomial Logistic Regression.

F Details of Experiments

We always set $K = 500$ for all of the first order methods. For RHG, TRHG and CG, we severally set $T = 100$ and $T = 50$ in toy example and real datasets. For CG method, $J = 200$ in toy example and $J = 20$ in real datasets. We let TRHG truncate at $T/2$ and set $T_z = 50, T_y = 25$ for BVFIM. BTRIPM severally uses 100 and 50 gradient descent steps to solve $\mathbf{z}^*(\mathbf{x})$ in toy example and real datasets. We use batch gradient descent in all experiments. We always set the step size for solving LL problem of different methods the same length. We scale the original data with a constant a for numerical stable consideration, where a varies for different datasets. We point out that the parameters τ, μ for BVFIM and BTRIPM are different because BVFIM and BTRIPM have different converging rates. Also, K varies for BTRIPM in different experiments depending on its convergence rate. For example, in "WIL" dataset, BTRIPM can converge in less than five steps, then we only need to set $K = 20$.

F.1 Toy example

$a = 1$. We set $\tau_k = 1/1.01^k, \mu_k = 10/1.04^k, s_1 = 0.1, s_2 = 0.01, \alpha = 0.01$ for BVFIM, where s_1, s_2, α are the step lengths as in Liu et al. [13]. We set $\tau_k = 1/1.02^k, \mu_k = 4/1.4^k, K = 100$ for BTRIPM.

F.2 MNIST dataset

$a = 0.0005$. We set $\tau_k = 1/1.03^k, \mu_k = 6/1.02^k, s_1 = 0.4, s_2 = 0.01, \alpha = 0.01$ for BVFIM1, $\tau_k = 1/1.01^k, \mu_k = 4/1.02^k, s_1 = 0.4, s_2 = 0.01, \alpha = 0.01$ for BVFIM2, $\tau_k = 1/1.3^k, \mu_k = 10/1.2^k, K = 50$ for BTRIPM1 and $\tau_k = 1/1.1^k, \mu_k = 8/1.2^k, K = 50$ for BTRIPM2.

F.3 FashionMNIST dataset

$a = 0.0004$. We set $\tau_k = 1/1.03^k, \mu_k = 5/1.02^k, s_1 = 0.1, s_2 = 0.01, \alpha = 0.01$ for BVFIM1 and BVFIM2, $\tau_k = 1/1.3^k, \mu_k = 6/1.2^k, K = 50$ for BTRIPM1 and BTRIPM2.

F.4 WIL dataset

$a = 0.001$. We set $\tau_k = 1/1.01^k, \mu_k = 4/1.04^k, s_1 = 0.1, s_2 = 0.01, \alpha = 0.01$ for BVFIM and $\tau_k = 1/1.02^k, \mu_k = 4/1.5^k, K = 20$ for BTRIPM.

788 **F.5 CM dataset**

789 $a = 0.0001$. We set $\tau_k = 1/1.01^k$, $\mu_k = 4/1.04^k$, $s_1 = 0.5$, $s_2 = 0.01$, $\alpha = 0.01$ for BVFIM and
 790 $\tau = 0.8/1.3^k$, $\mu_k = 4/1.2^k$, $K = 25$ for BTRIPM.

791 **F.6 20 Newsgroup**

792 $a = 1$. We set $\tau_k = 1/1.02^k$, $\mu_k = 10/1.02^k$, $s_1 = 0.4$, $s_2 = 0.01$, $\alpha = 0.01$ for BVFIM and
 793 $\tau_k = 1/1.2^k$, $\mu_k = 10/1.2^k$, $K = 50$ for BTRIPM.

794 **G Validation of Our Assumptions on the Toy Example (12)**

795 **Assumption 3.1:** We consider $\mu_k < 1$ and (x_k, y_k) locally minimizes the toy example. By simple
 796 analysis we know

$$\sin(x_k + y_k) = -1 + \mu_k$$

797 must hold. Further $x_k + y_k = \arcsin(\mu_k - 1) + 2n\pi$ when $n \leq 0$ and $x_k + y_k = -\arcsin(\mu_k -$
 798 $1) + (2n - 1)\pi$ when $n > 0$. Here n is a constant integer. It suffices to test our assumption for $n \leq 0$.
 799 It is easy to know

$$x_k = y_k = \frac{\arcsin(\mu_k - 1)}{2} + n\pi.$$

800 Then $\{(x_k, y_k)\}$ is bounded and

$$(x_k, y_k) \rightarrow (\bar{x}, \bar{y}) = -\frac{\pi}{4} + n\pi$$

801 as $\mu_k \rightarrow 0$. Absolutely $\partial_{\mathbf{y}\mathbf{y}}^2 f(\bar{x}, \bar{y}) = -\sin(\bar{x}, \bar{y}) = 1 > 0$. Then we have proved (1) and (2). For
 802 (3) we can choose $z^*(x) = -\frac{\pi}{2} + 2n\pi - x$ in the neighborhood of \bar{x} , then $z^*(\bar{x}) = \bar{y}$. Since in this
 803 example $\mathcal{X} = \mathbb{R}$, the linear independence assumption automatically holds. Therefore, Assumption
 804 3.1 holds for this example.

805 To deeply investigate (3) in Assumption 3.1, we establish the following Lemma. This Lemma reveals
 806 that (3) in Assumption 3.1 is weaker than the global strong convexity assumption of $f(\mathbf{x}, \mathbf{y})$ w.r.t. \mathbf{y} .

807 **Lemma G.1.** *If for all $\mathbf{x} \in \mathcal{X}$ $f(\mathbf{x}, \mathbf{y})$ is strongly convex w.r.t. \mathbf{y} and $\operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is nonempty,*
 808 *then $\mathbf{z}^*(\mathbf{x})$ is unique and is a continuous function of \mathbf{x} . Moreover, $\mathbf{z}^*(\mathbf{x})$ is differentiable.*

809 *Proof.* Since $f(\mathbf{x}, \mathbf{y})$ is strongly convex w.r.t. \mathbf{y} and $\operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is nonempty, we conclude that
 810 $\operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is a singleton. Then $\mathbf{z}^*(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is unique and $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{z}^*(\mathbf{x})) = \mathbf{0}$.
 811 Note that \mathbf{y} satisfying $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}$ is also unique, we claim that $\mathbf{y} = \mathbf{z}^*(\mathbf{x})$ is equivalent to
 812 $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}$. Then by implicit differentiation theorem, we deduce that $\mathbf{z}^*(\mathbf{x})$ is a differentiable
 813 function of \mathbf{x} in a neighborhood of \mathbf{x} because $\partial_{\mathbf{y}\mathbf{y}}$ is invertible due to the strong convexity of f w.r.t
 814 \mathbf{y} . Since \mathbf{x} is an arbitrary vector in \mathcal{X} , $\mathbf{z}^*(\mathbf{x})$ is a differentiable function of \mathbf{x} in \mathcal{X} . \square