# Binary Optimization via Mathematical Programming
# with Equilibrium Constraints

Ganzhao Yuan [*]　　　Bernard Ghanem [†]

## Abstract

*Binary optimization is a central problem in mathematical optimization and its applications are abundant. To solve this problem, we propose a new class of continuous optimization techniques which is based on Mathematical Programming with Equilibrium Constraints (MPECs). We first reformulate the binary program as an equivalent augmented biconvex optimization problem with a bilinear equality constraint, then we propose two penalization/regularization methods (exact penalty and alternating direction) to solve it. The resulting algorithms seek desirable solutions to the original problem via solving a sequence of linear programming convex relaxation subproblems. In addition, we prove that both the penalty function and augmented Lagrangian function, induced by adding the complementarity constraint to the objectives, are exact, i.e., they have the same local and global minima with those of the original binary program when the penalty parameter is over some threshold. The convergence of both algorithms can be guaranteed, since they essentially reduce to block coordinate descent in the literature. Finally, we demonstrate the effectiveness and versatility of our methods on several important problems, including graph bisection, constrained image segmentation, dense subgraph discovery, modularity clustering and Markov random fields. Extensive experiments show that our methods outperform existing popular techniques, such as iterative hard thresholding, linear programming relaxation and semidefinite programming relaxation.*

***Keywords:*** *Binary Optimization, Convergence Analysis, MPECs, Exact Penalty Method, Alternating Direction Method, Graph Bisection, Constrained Image Segmentation, Dense Subgraph Discovery, Modularity Clustering, Markov Random Fields.*

## 1. Introduction

In this paper, we mainly focus on the following binary optimization problem:

$$\min_{\mathbf{x}} \ f(\mathbf{x}), \ s.t. \ \mathbf{x} \in \{-1, 1\}^n, \ \mathbf{x} \in \Omega \qquad (1)$$

where the objective function $f : \mathbb{R}^n \to \mathbb{R}$ is convex (but not necessarily smooth) on some convex set $\Omega$, and the non-convexity of Eq (1) is only caused by the binary constraints. In addition, we assume $\{-1, 1\}^n \cap \Omega \neq \emptyset$.

The optimization in Eq (1) describes many applications of interest in both computer vision and machine learning, including graph bisection [24, 35], image (co-)segmentation [51, 35, 33], Markov random fields [8], permutation problem [21], graph matching [17, 64, 56], binary matrix completion [18, 29], hashing coding [39, 59], image registration [60], multimodal feature learning [54], multi-target tracking [52], visual alignment [53], and social network analysis (e.g. subgraphs discovery [68, 2], biclustering [1], planted k-disjoint-clique discover [4], planted clique and biclique discovery [3], community discovery [27, 14]), etc.

The binary optimization problem is difficult to solve, since it is NP-hard. One type of methods to solve this problem is continuous in nature. The simple way is to relax the binary constraint with Linear Programming (LP) relaxation constraints $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$ and round the entries of the resulting continuous solution to the nearest integer at the end. However, not only may this solution not be optimal, it may not even be feasible and violate some constraint. Another type of optimization focuses on the cutting-plane and branch-and-cut method. The cutting plane method solves the LP relaxation and then adds linear constraints that drive the solution towards integers. The branch-and-cut method partially develops a binary tree and iteratively cuts out the enumeration that is worse than current lower bound, while the lower bound can be found using convex relaxation, Lagrangian duality, or Lipschitz continuity. However, this class of method ends up solving all $2^n$ convex subproblems in the worst case. Our algorithm aligns with the first research direction. It relies on solving a convex LP relaxation subproblem iteratively, but it provably terminates in finite iterations.

---

[*]King Abdullah University of Science & Technology (KAUST). Email: yuanganzhao@gmail.com.

[†]King Abdullah University of Science & Technology (KAUST). Email: bernard.ghanem@kaust.edu.sa.

Table 1: Existing continuous methods for binary optimization.

| | Method and Reference | Description |
|---|---|---|
| **Relaxed Approximation** | spectral relaxation [16, 51] | $\{-1,+1\}^n \approx \{\mathbf{x} \mid \|\mathbf{x}\|_2^2 = n\}$ |
| | linear programming relaxation [36, 29] | $\{-1,+1\}^n \approx \{\mathbf{x} \mid -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}\}$ |
| | SDP relaxation [60, 61, 35] | $\{0,+1\}^n \approx \{\mathbf{x} \mid \mathbf{X} \succeq \mathbf{x}\mathbf{x}^T, diag(\mathbf{X}) = \mathbf{x}\}$ |
| | | $\{-1,+1\}^n \approx \{\mathbf{x} \mid \mathbf{X} \succeq \mathbf{x}\mathbf{x}^T, \ diag(\mathbf{X}) = \mathbf{1}\}$ |
| | doubly positive relaxation [31, 61] | $\{0,+1\}^n \approx \{\mathbf{x} \mid \mathbf{X} \succeq \mathbf{x}\mathbf{x}^T, diag(\mathbf{X}) = \mathbf{x}, \ \mathbf{x} \geq \mathbf{0}, \ \mathbf{X} \geq \mathbf{0}\}$ |
| | completely positive relaxation [12, 11] | $\{0,+1\}^n \approx \{\mathbf{x} \mid \mathbf{X} \succeq \mathbf{x}\mathbf{x}^T, diag(\mathbf{X}) = \mathbf{x}, \ \mathbf{x} \geq \mathbf{0}, \ \mathbf{X} \text{ is CP}\}$ |
| | SOCP relaxation [37, 23] | $\{-1,+1\}^n \approx \{\mathbf{x} \mid \langle \mathbf{X} - \mathbf{x}\mathbf{x}^T, \mathbf{L}\mathbf{L}^T \rangle \geq 0, \ diag(\mathbf{X}) = \mathbf{1}\}, \ \forall \mathbf{L}$ |
| **Equivalent Optimization** | iterative hard thresholding [68, 5] | $\min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}'\|_2^2, \ s.t. \ \mathbf{x} \in \{-1,+1\}^n$ |
| | $\ell_0$ norm reformulation [40, 67] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid \|\mathbf{x}+\mathbf{1}\|_0 + \|\mathbf{x}-\mathbf{1}\|_0 \leq n\}$ |
| | piecewise separable reformulation [70] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid (\mathbf{1}+\mathbf{x}) \odot (\mathbf{1}-\mathbf{x}) = \mathbf{0}\}$ |
| | $\ell_2$ box non-separable reformulation [48, 42] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ \|\mathbf{x}\|_2^2 = n\}$ |
| | $\ell_p$ box non-separable reformulation [63] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ \|\mathbf{x}\|_p^p = n, \ 0 < p < \infty\}$ |
| | $\ell_\infty$ box separable MPEC $\quad$ [This paper] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid \mathbf{x} \odot \mathbf{v} = \mathbf{1}, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_\infty \leq 1, \ \forall \mathbf{v}\}$ |
| | $\ell_2$ box separable MPEC $\quad$ [This paper] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid \mathbf{x} \odot \mathbf{v} = 1, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_2^2 \leq n, \ \forall \mathbf{v}\}$ |
| | $\ell_\infty$ box non-separable MPEC $\quad$ [This paper] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid \langle \mathbf{x}, \mathbf{v} \rangle = n, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_\infty \leq 1, \ \forall \mathbf{v}\}$ |
| | $\ell_2$ box non-separable MPEC $\quad$ [This paper] | $\{-1,+1\}^n \Leftrightarrow \{\mathbf{x} \mid \langle \mathbf{x}, \mathbf{v} \rangle = n, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_2^2 \leq n, \ \forall \mathbf{v}\}$ |

In non-convex optimization, good initialization is very important to the quality of the solution. Motivated by this, several papers design smart initialization strategies and establish optimality qualification of the solutions for non-convex problems. For example, the work of [69] considers a multi-stage convex optimization algorithm to refine the global solution by the initial convex method; the work of [13] starts with a careful initialization obtained by a spectral method and improves this estimate by gradient descent; the work of [32] uses the top-$k$ singular vectors of the matrix as initialization and provides theoretical guarantees for biconvex alternating minimization algorithm. The proposed method also uses a similar initialization strategy since it reduces to convex LP relaxation in the first iteration.

The contributions of this paper are three-fold. **(a)** We reformulate the binary program as an equivalent augmented optimization problem with a bilinear equality constraint via a variational characterization of the binary constraint. Then, we propose two penalization/regularization methods (exact penalty and alternating direction) to solve it. The resulting algorithms seek desirable solutions to the original binary program. **(b)** We prove that both the penalty function and augmented Lagrangian function, induced by adding the complementarity constraint to the objectives, are exact, i.e., the set of their globally optimal solutions coincide with that of Eq (1) when the penalty parameter is over some threshold. Thus, the convergence of both algorithms can be guaranteed since they reduce to block coordinate descent in the literature [57, 7]. This is the first attempt to solve general non-smooth binary program with guaranteed convergence. **(c)** We provide numerical comparisons with state-of-the-

art techniques, such as iterative hard thresholding [68], linear programming relaxation [36, 37] and semidefinite programming relaxation [60] on a variety of concrete computer vision and machine learning problems. Extensive experiments have demonstrated the effectiveness of our proposed methods.

This paper is organized as follows. Section 2 provides a brief description of the related work. Section 3 presents our MPEC-based optimization framework. Section 4 discusses some features of our methods. Section 5 summarizes the experimental results. Finally, Section 6 concludes this paper. Throughout this paper, we use lowercase and uppercase boldfaced letters to denote real vectors and matrices respectively. The Euclidean inner product between $\mathbf{x}$ and $\mathbf{y}$ is denoted by $\langle \mathbf{x}, \mathbf{y} \rangle$ or $\mathbf{x}^T\mathbf{y}$. We use $\mathbf{I}_n$ to denote an identity matrix of size $n$, where sometimes the subscript is dropped when $n$ is known from the context. $\mathbf{X} \succeq \mathbf{0}$ means that matrix $\mathbf{X}$ is positive semi-definite. Finally, sign is a signum function with $\text{sign}(0) = \pm 1$.

## 2. Related Work

This paper proposes a new continuous method for binary optimization. We briefly review existing related work in this research direction in the literature (see Table 1).

There are generally two types of methods in the literature. One is the relaxed approximation method. Spectral relaxation [16, 46, 51, 38] replaces the binary constraint with a spherical constraint and solves the problem using eigen decomposition. Despite its computational merits, it is difficult to generalize to handle linear or nonlinear con-

straints. Linear programming relaxation [36, 37, 10] transforms the NP-hard optimization problem into a convex box-constrained optimization problem, which can be solved by well-established optimization methods and software [25]. Semi-Definite Programming (SDP) relaxation [37, 31] uses a lifting technique $\mathbf{X} = \mathbf{x}\mathbf{x}^T$ and relaxes to a convex conic $\mathbf{X} \succeq \mathbf{x}\mathbf{x}^T$ [1] to handle the binary constraint. Combining this with a unit-ball randomized rounding algorithm, the work of [24] proves that at least a factor of 87.8% to the global optimal solution can be achieved for the graph bisection problem. Since the original paper of [24], SDP has been applied to develop numerous approximation algorithms for NP-hard problems. As more constraints lead to tighter bounds for the objective, doubly positive relaxation considers constraining both the eigenvalues and the elements of the SDP solution to be nonnegative, leading to better solution than canonical SDP method. In addition, Completely Positive (CP) relaxation [12, 11] further constrains the entries of the factorization of the solution $\mathbf{X} = \mathbf{L}\mathbf{L}^T$ to be nonnegative $\mathbf{L} \geq \mathbf{0}$. It can be solved by tackling its associated dual co-positive program, which is related to the study of indefinite optimization and sum-of-squares optimization in the literature. Second-Order Cone Programming (SOCP) relaxes the SDP conic into the nonnegative orthant [37] using the fact that $\langle \mathbf{X} - \mathbf{x}\mathbf{x}^T, \mathbf{L}\mathbf{L}^T \rangle \geq 0, \ \forall \mathbf{L}$, resulting in tighter bound than the LP method, but looser bound than that of the SDP method. Therefore it can be viewed as a balance between efficiency and efficacy.

Another types of methods for binary optimization relates to equivalent optimization. The iterative hard thresholding method directly handles the non-convex constraint via projection and it has been widely used due to its simplicity and efficiency [68, 5, 40]. However, this method is often observed to obtain sub-optimal accuracy and it is not directly applicable, when the objective function is non-smooth. Binary optimization can be reformulated as an $\ell_0$ norm semi-continuous optimization problem [2]. Thus, existing $\ell_0$ norm sparsity constrained optimization techniques such as quadratic penalty decomposition method [40, 5] and multi-stage convex optimization method [69, 67] can be applied. A piecewise separable reformulation has been considered in [70], which can exploit existing smooth optimization techniques. A continuous $\ell_2$ box non-separable reformulation [3] has been used in the literature [48, 34, 47]. A second-order interior point method [42, 19] has been developed to solve the continuous reformulation optimization problem. A continuous $\ell_p$ box non-separable reformulation has been used in [63], where an interesting geometric illustration of $\ell_p$-box intersection has been shown [4]. In addition, they infuse this equivalence into the optimization framework of Alternating Direction Method of Multipliers (ADMM). However, their guarantee of convergence is weak. In this paper, to tackle the binary optimization problem, we propose a new framework that is based on Mathematical Programming with Equilibrium Constraints (MPECs) (refer to the proposed MPEC reformulations in Table 1 [5]). Our resulting algorithms are theoretically convergent and empirically effective.

Mathematical programs with equilibrium constraints [6] are optimization problems where the constraints include complementarities or variational inequalities. They are difficult to deal with because their feasible region may not necessarily be convex or even connected. Motivated by recent development of MPECs for non-convex optimization [65, 66, 67, 6, 41], we consider continuous $\ell_2$ box non-separable MPEC of the binary optimization problem. In our forthcoming algorithm design, we mainly focus on this specific formulation.

## 3. Proposed Optimization Algorithms

This section presents our MPEC-based optimization algorithms. We first propose an equivalent reformulation of binary optimization program, and then we consider two algorithms (exact penalty method and alternating direction method) to solve it.

### 3.1. Equivalent Reformulation

First of all, we present a variational reformulation of the binary constraint [7].

**Lemma 1.** $\boldsymbol{\ell_2}$ **box non-separable MPEC.** *We define* $\Theta \triangleq \{(\mathbf{x}, \mathbf{v}) \mid \mathbf{x}^T\mathbf{v} = n, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_2^2 \leq n\}$. *Assume that* $(\mathbf{x}, \mathbf{v}) \in \Theta$, *then we have* $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$, *and* $\mathbf{x} = \mathbf{v}$.

Using Lemma 1, we can rewrite Eq (1) in an equivalent form as follows.

$$\min_{-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ \|\mathbf{v}\|_2^2 \leq n} f(\mathbf{x}), \ s.t. \ \mathbf{x}^T\mathbf{v} = n, \ \mathbf{x} \in \Omega \qquad (2)$$

We remark that $\mathbf{x}^T\mathbf{v} = n$ is called complementarity (or equilibrium) constraint in the literature [41, 49, 55] and it always holds that $\mathbf{x}^T\mathbf{v} \leq \|\mathbf{x}\|_\infty\|\mathbf{v}\|_1 \leq \sqrt{n}\|\mathbf{v}\|_2 \leq n$ for any feasible $\mathbf{x}$ and $\mathbf{v}$.

---

[1]Using Schur complement lemma: $\forall \mathbf{A}, \mathbf{B}, \mathbf{C}, \ \mathbf{A} \succeq \mathbf{B}^T\mathbf{C}^\dagger\mathbf{B}, \ \mathbf{C} \succeq \mathbf{0} \Leftrightarrow \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix} \succeq \mathbf{0}$, one can rewrite $\mathbf{X} \succeq \mathbf{x}\mathbf{x}^T$ as $\begin{pmatrix} \mathbf{X} & \mathbf{x} \\ \mathbf{x}^T & 1 \end{pmatrix} \succeq \mathbf{0}$.

[2]One can rewrite the $\ell_0$ norm constraint $\|\mathbf{x} + \mathbf{1}\|_0 + \|\mathbf{x} - \mathbf{1}\|_0 \leq n$ in a compact matrix form $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_0 \leq n$ with $\mathbf{A} = [\mathbf{I}_n \mid \mathbf{I}_n]^T \in \mathbb{R}^{2n \times n}$, $\mathbf{b} = [\mathbf{1}^T \mid -\mathbf{1}^T]^T \in \mathbb{R}^{2n}$.

[3]They replace $\mathbf{x} \in \{0, 1\}^n$ with $\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}, \ \mathbf{x}^T(\mathbf{1} - \mathbf{x}) = 0$. We extend this strategy to replace $\{-1, +1\}^n$ with $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ (\mathbf{1} + \mathbf{x})^T(\mathbf{1} - \mathbf{x}) = 0$ which can be simplified to $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ \|\mathbf{x}\|_2^2 = n$. Appendix A provides a simple proof.

[4]Here we adapt their $\{0, 1\}$ formulation to our $\{-1, +1\}$ formulation.

[5]Appendix A provides the proofs of all the equivalent reformations.

[6]In fact, we focus on mathematical programs with complementarity constraints (MPCC), a subclass of MPECs where the original discrete optimization problem can be formulated as a complementarity problem and therefore as a nonlinear program. Here we use the term MPECs for the purpose of generality and historic conventions.

[7]All proofs can be found in the Appendix.

## 3.2. Exact Penalty Method

We now present our exact penalty method for solving the optimization problem in Eq (2). It is worthwhile to point out that there are many studies on exact penalty for MPECs (refer to [41, 30, 49, 55, 67] for examples), but they do not afford the exactness of our penalty problem. In an exact penalty method, we penalize the complementary error directly by a penalty function. The resulting objective $\mathcal{J} : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ is defined in Eq (3), where $\rho$ is the penalty parameter that is iteratively increased to enforce the constraints.

$$\mathcal{J}_\rho(\mathbf{x}, \mathbf{v}) = f(\mathbf{x}) + \rho(n - \mathbf{x}^T \mathbf{v})$$
$$s.t. \ -\mathbf{1} \le \mathbf{x} \le \mathbf{1}, \ \|\mathbf{v}\|_2^2 \le n, \ \mathbf{x} \in \Omega \quad (3)$$

In each iteration, we minimize over $\mathbf{x}$ and $\mathbf{v}$ alternatingly [57, 7], while fixing the parameter $\rho$. We summarize our exact penalty method in Algorithm 1. The parameter $T$ is the number of inner iterations for solving the biconvex problem. We make the following observations about the algorithm.

**(a) Initialization**. We initialize $\mathbf{v}^0$ to $\mathbf{0}$. This is for the sake of finding a reasonable local minimum in the first iteration, as it reduces to LP convex relaxation [36] for the binary optimization problem.

**(b) Exact property**. One remarkable feature of our method is the boundedness of the penalty parameter $\rho$ (see Theorem 1). Therefore, we terminate the optimization when the threshold is reached (see Eq (8)). This distinguishes it from the quadratic penalty method [40], where the penalty may become arbitrarily large for non-convex problems.

**(c) $\mathbf{v}$-Subproblem**. Variable $\mathbf{v}$ in Eq (7) is updated by solving the following convex problem:

$$\mathbf{v}^{t+1} = \arg\min \ \langle \mathbf{v}, -\mathbf{x}^{t+1} \rangle, \ s.t. \ \|\mathbf{v}\|_2^2 \le n \quad (4)$$

When $\mathbf{x}^{t+1} = 0$, any feasible solution is also an optimal solution; when $\mathbf{x}^{t+1} \ne 0$, the optimal solution will be achieved in the boundary with $\|\mathbf{v}\|_2^2 = n$ and Eq (4) is equivalent to solving: $\min_{\|\mathbf{v}\|_2^2 = n} \frac{1}{2}\|\mathbf{v}\|_2^2 - \langle \mathbf{v}, \mathbf{x}^{t+1} \rangle$. Thus, we have the following optimal solution for $\mathbf{v}$:

$$\mathbf{v}^{t+1} = \begin{cases} \sqrt{n} \cdot \mathbf{x}^{t+1}/\|\mathbf{x}^{t+1}\|_2, & \mathbf{x}^{t+1} \ne 0; \\ \text{any } \mathbf{v} \text{ with } \|\mathbf{v}\|_2^2 \le n, & \text{otherwise.} \end{cases} \quad (5)$$

**(d) $\mathbf{x}$-Subproblem**. Variable $\mathbf{x}$ in Eq (6) is updated by solving box constrained convex problem which has no closed-form solution. However, it can be solved using Nesterov's proximal gradient method [43] or classical/linearized ADM [26].

**Theoretical Analysis.** In Theorem 1, we show that when the penalty parameter $\rho$ is larger than some threshold, the biconvex objective function in Eq (3) is equivalent to the original constrained MPEC problem in Eq (2). This essentially implies the theoretical convergence of the algorithm

**Algorithm 1** MPEC-EPM: An Exact Penalty Method for Solving MPEC Problem (2)

---

(S.0) Set $t = 0$, $\mathbf{x}^0 = \mathbf{v}^0 = \mathbf{0}$, $\rho > 0$, $\sigma > 1$.
(S.1) Solve the following $\mathbf{x}$-subproblem [primal step]:

$$\mathbf{x}^{t+1} = \arg\min_{\mathbf{x}} \ \mathcal{J}(\mathbf{x}, \mathbf{v}^t), \ s.t. \ -\mathbf{1} \le \mathbf{x} \le \mathbf{1}, \ \mathbf{x} \in \Omega \quad (6)$$

(S.2) Solve the following $\mathbf{v}$-subproblem [dual step]:

$$\mathbf{v}^{t+1} = \arg\min_{\mathbf{v}} \ \mathcal{J}(\mathbf{x}^{t+1}, \mathbf{v}), \ s.t. \ \|\mathbf{v}\|_2^2 \le n \quad (7)$$

(S.3) Update the penalty in every $T$ iterations:

$$\rho \Leftarrow \min(2L, \rho \times \sigma) \quad (8)$$

(S.4) Set $t := t + 1$ and then go to Step (S.1)

---

since it reduces to block coordinate descent in the literature [8]. Moreover, Theorem 2 characterizes the convergence rate and asymptotic monotone property the algorithm [9].

**Theorem 1.** *Exactness of the Penalty Function. Assume that $f(\cdot)$ is a L-Lipschitz continuous convex function on $-\mathbf{1} \le \mathbf{x} \le \mathbf{1}$. When $\rho > 2L$, the biconvex optimization $\min_{\mathbf{x}, \mathbf{v}} \mathcal{J}_\rho(\mathbf{x}, \mathbf{v}), \ s.t. \ -\mathbf{1} \le \mathbf{x} \le \mathbf{1}, \|\mathbf{v}\|_2^2 \le n, \ \mathbf{x} \in \Omega$ in Eq (3) has the same local and global minima with the original problem in Eq (2).*

**Theorem 2.** *Convergence Rate and Asymptotic Monotone Property of Algorithm 1. Assume that $f(\cdot)$ is a L-Lipschitz continuous convex function on $-\mathbf{1} \le \mathbf{x} \le \mathbf{1}$. Algorithm 1 will converge to the first-order KKT point in at most $\lceil (\ln(L\sqrt{2n}) - \ln(\epsilon\rho^0))/\ln\sigma \rceil$ outer iterations [10] with the accuracy at least $n - \mathbf{x}^T \mathbf{v} \le \epsilon$. Moreover, after $\langle \mathbf{x}, \mathbf{v} \rangle = n$ is obtained, the sequence of $\{f(\mathbf{x}^t)\}$ generated by Algorithm 1 is monotonically non-increasing.*

## 3.3. Alternating Direction Method

This section presents an alternating direction method (of multipliers) (ADM or ADMM) for solving Eq (2). This is mainly motivated by the recent popularity of ADM in the optimization literature [26, 61, 65, 66].

We first form the augmented Lagrangian $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}$ in Eq (9) as follows:

$$\mathcal{L}(\mathbf{x}, \mathbf{v}, \rho) \triangleq f(\mathbf{x}) + \rho(n - \mathbf{x}^T \mathbf{v}) + \frac{\alpha}{2}(n - \mathbf{x}^T \mathbf{v})^2$$
$$s.t. \ -\mathbf{1} \le \mathbf{x} \le \mathbf{1}, \ \|\mathbf{v}\|_2^2 \le n, \ \mathbf{x} \in \Omega \quad (9)$$

---

[8]Specifically, using Tseng's convergence results of block coordinate descent for non-differentiable minimization [57], one can guarantee that every clustering point of Algorithm 2 is also a stationary point. In addition, stronger convergence results [7, 67] can be obtained by combining a proximal strategy and Kurdyka-Łojasiewicz inequality assumption on $\mathcal{J}(\cdot)$.

[9]Appendix B provides all the proofs of the convergence theorems.

[10]Every time we increase $\rho$, we call it one outer iteration.

where $\rho$ is the Lagrange multiplier associated with the complementarity constraint $n - \langle \mathbf{x}, \mathbf{v} \rangle = 0$, and $\alpha > 0$ is the penalty parameter. Interestingly, we find that the augmented Lagrangian function can be viewed as adding an elastic net regularization [71] on the complementarity error. We detail the ADM iteration steps for Eq (9) in Algorithm 2, which has the following properties.

**(a) Initialization**. We set $\mathbf{v}^0 = \mathbf{0}$ and $\rho^0 = 0$. This finds a reasonable local minimum in the first iteration, as it reduces to LP relaxation for the $\mathbf{x}$-subproblem.

**(b) Monotone property**. For any feasible solution $\mathbf{x}$ and $\mathbf{v}$ in Eq (9), it holds that $n - \mathbf{x}^T \mathbf{v} \geq 0$. Using the fact that $\alpha^t > 0$ and due to the $\rho^t$ update rule, $\rho^t$ is monotone increasing.

**(c) $\mathbf{v}$-Subproblem**. Variable $\mathbf{v}$ in Eq (11) is updated by solving the following problem:

$$\mathbf{v}^{t+1} = \arg\min_{\mathbf{v}} \frac{1}{2} \mathbf{v}^T \mathbf{a} \mathbf{a}^T \mathbf{v} + \langle \mathbf{v}, \mathbf{b} \rangle, \; s.t. \; \|\mathbf{v}\|_2 \leq \sqrt{n}$$

where $\mathbf{a} \triangleq \mathbf{x}^{t+1}/\sqrt{\alpha}$ and $\mathbf{b} \triangleq -(\rho + \alpha n) \cdot \mathbf{x}^{t+1}$. This problem is also known as constrained eigenvalue problem in the literature [22], which has efficient solution. Since the quadratic matrix is of rank one, this problem has nearly closed-form solution (please refer to Appendix C).

**(d) $\mathbf{x}$-Subproblem**. Variable $\mathbf{x}$ in Eq (9) is updated by solving a box constrained optimization problem. Similar to Eq (6), it can be solved using Nesterov's proximal gradient method or classical/linearized ADM.

**Theoretical Analysis.** Our MPEC-ADM has excellent convergence property. In Theorem 3, we show that when the multiplier $\rho$ is larger than some threshold, the biconvex objective function in Eq (9) is equivalent to the original constrained MPEC problem in Eq (2). The theoretical convergence of the algorithm directly follows as it is similar to our previous analysis for the exact penalty method.

**Theorem 3.** *Exactness of the augmented Lagrangian Function. Assume that $f(\cdot)$ is a L-Lipschitz continuous convex function on $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. When $\rho > 2L$, the biconvex optimization problem $\min_{\mathbf{x}, \mathbf{v}} \mathcal{L}(\mathbf{x}, \mathbf{v}, \rho), \; s.t. \; -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \; \|\mathbf{v}\|_2^2 \leq n, \; \mathbf{x} \in \Omega$ in Eq (9) has the same local and global minima with the original problem in Eq (2).*

## 4. Discussions

**MPEC-EPM vs. MPEC-ADM.** The proposed MPEC-EPM and MPEC-ADM have their own advantages. (a) MPEC-EPM is more simple and elegant and it can directly use existing LP relaxation optimization solver. (b) MPEC-EPM may be less adaptive since the penalty parameter $\rho$ is monolithically increased until a threshold is achieved. In comparison, MPEC-ADM is more adaptive, since a constant penalty also guarantees monotonically non-decreasing multipliers and convergence.

---

**Algorithm 2** MPEC-ADM: An Alternating Direction Method for Solving MPEC Problem (2)

---

(S.0) Set $t = 0$, $\mathbf{x}^0 = \mathbf{v}^0 = \mathbf{0}$, $\rho^0 = 0$, $\alpha > 0$, $\sigma > 1$.

(S.1) Solve the following $\mathbf{x}$-subproblem [primal step]:

$$\mathbf{x}^{t+1} = \arg\min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{v}^t, \rho^t), \; s.t. \; \mathbf{x} \in [-1, +1]^n \cap \Omega \quad (10)$$

(S.2) Solve the following $\mathbf{v}$-subproblem [dual step]:

$$\mathbf{v}^{t+1} = \arg\min_{\mathbf{v}} \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{v}, \rho^t) \; s.t. \; \|\mathbf{v}\|_2^2 \leq n \quad (11)$$

(S.3) Update the Lagrange multiplier:

$$\rho^{t+1} = \rho^t + \alpha(n - \langle \mathbf{x}^{t+1}, \mathbf{v}^{t+1} \rangle) \quad (12)$$

(S.4) Update the penalty in every $T$ iterations (if necessary):

$$\alpha \Leftarrow \alpha \times \sigma \quad (13)$$

(S.5) Set $t := t + 1$ and then go to Step (S.1).

---

**Merits of our methods.** There are several merits behind our MPEC-based penalization/regularization methods. (a) They exhibit strong convergence guarantees since they essentially reduce to block coordinate descent in the literature [57, 7]. (b) They seek desirable solutions since the LP convex relaxation methods in the first iteration provide good initializations. (c) They are efficient since they are amenable to the use of existing convex methods to solve the sub-problems. (d) They have a monotone/greedy property due to the complimentary constraints brought on by MPECs. We penalize the complimentary error and ensure that the error is decreasing in every iteration, leading to binary solutions.

**Extensions to Zero-One Constraint and Orthogonality Constraint.** For convenience, we define $\Delta = \{0, 1\}^n$ and $\mathbb{O} = \{\mathbf{X} \in \mathbb{R}^{n \times r} \mid \mathbf{X}^T \mathbf{X} = \mathbf{I}_r \; (n \geq r)\}$. Noticing that $\mathbf{y} = (\mathbf{x} + \mathbf{1})/2 \in \{0, 1\}^n$, we can extend $\ell_\infty$ box non-separable MPEC [11] and $\ell_2$ box non-separable MPEC [12] to handle zero-one binary optimization. Moreover, observing that binary constraint $\mathbf{x} \in \{-1, +1\}^n \Leftrightarrow |\mathbf{x}| = \mathbf{1}$ is analogous to orthogonality constraint since $\mathbf{X} \in \mathbb{O} \Leftrightarrow \sigma(\mathbf{X}) = \mathbf{1}$, we can extend our $\ell_\infty$ box non-separable MPEC [13] and $\ell_2$ box non-separable MPEC [14] to the optimization problem with orthogonality constraint [62, 15].

---

[11] $\Delta \Leftrightarrow \{\mathbf{x} \mid \mathbf{0} \leq \mathbf{x} \leq \mathbf{1}, \; \mathbf{0} \leq \mathbf{v} \leq \mathbf{1}, \; \langle 2\mathbf{x} - \mathbf{1}, 2\mathbf{v} - \mathbf{1} \rangle = n, \; \forall \mathbf{v}\}$
[12] $\Delta \Leftrightarrow \{\mathbf{x} \mid \mathbf{0} \leq \mathbf{x} \leq \mathbf{1}, \; \|2\mathbf{v} - \mathbf{1}\|_2^2 \leq n, \; \langle 2\mathbf{x} - \mathbf{1}, 2\mathbf{v} - \mathbf{1} \rangle = n, \; \forall \mathbf{v}\}$
[13] $\mathbb{O} \Leftrightarrow \{\mathbf{X} \mid \mathbf{X}^T \mathbf{X} \preceq \mathbf{I}, \; \mathbf{V}^T \mathbf{V} \preceq \mathbf{I}, \; \langle \mathbf{X}, \mathbf{V} \rangle = r, \; \forall \mathbf{V}\}$
[14] $\mathbb{O} \Leftrightarrow \{\mathbf{X} \mid \mathbf{X}^T \mathbf{X} \preceq \mathbf{I}_r, \; \|\mathbf{V}\|_F^2 \leq r, \; \langle \mathbf{X}, \mathbf{V} \rangle = r, \; \forall \mathbf{V}\}$

| (a) LP | (b) NCUT | (c) RCUT | (d) SDP |
| $f = 473.646$ | $f = 230.049$ | $f = 548.964$ | $f = 194.664$ |

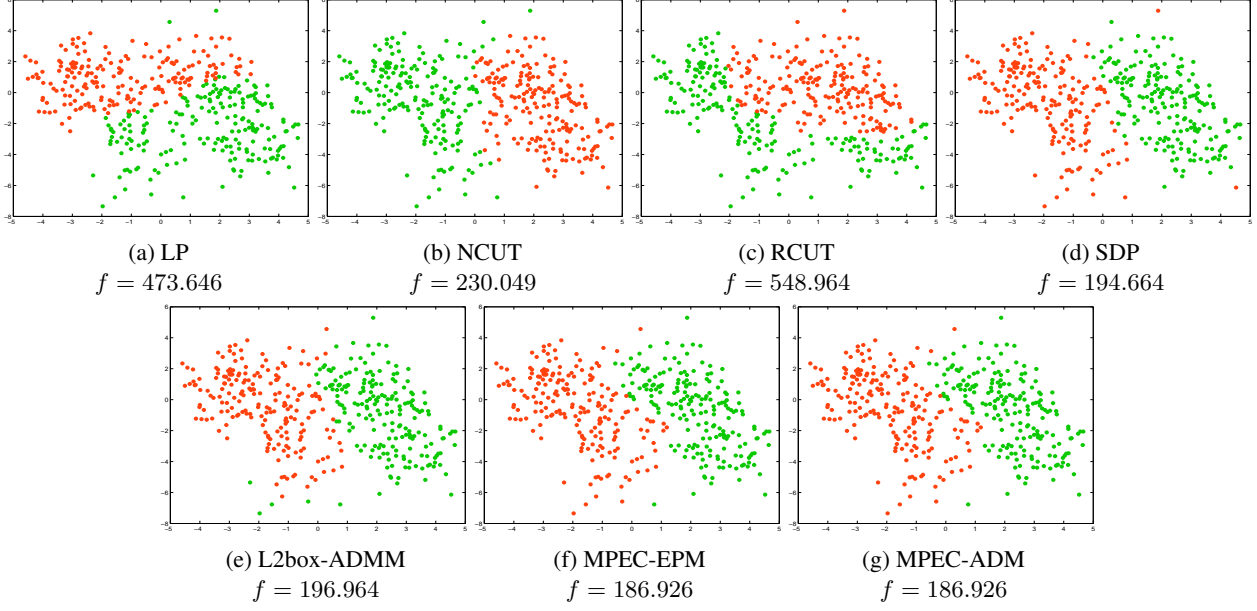| (e) L2box-ADMM | (f) MPEC-EPM | (g) MPEC-ADM |
| $f = 196.964$ | $f = 186.926$ | $f = 186.926$ |

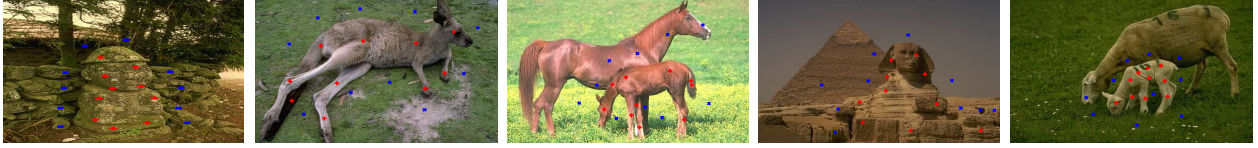Figure 1: Graph bisection on the '4gauss' data set.



Figure 2: Images in our constrained image segmentation experiments. 10 foreground pixels and 10 background pixels are annotated by red and blue markers respectively.

## 5. Experimental Validation

In this section, we demonstrate the effectiveness of our algorithms (MPEC-EPM and MPEC-ADM) on 5 binary optimization tasks, namely graph bisection, constrained image segmentation, dense subgraph discovery, modularity clustering and Markov random fields. All codes are implemented in MATLAB using a 3.20GHz CPU and 8GB RAM.

### 5.1. Graph Bisection

Graph bisection aims at separating the vertices of a weighted undirected graph into two disjoint sets with minimal cut edges with equal size. Mathematically, it can be formulated as the following optimization problem [35, 60]:

$$\min_{\mathbf{x} \in \{-1, +1\}^n} \mathbf{x}^T \mathbf{L} \mathbf{x}, \ s.t. \ \mathbf{x}^T \mathbf{1} = 0$$

where $\mathbf{L} = \mathbf{D} - \mathbf{W} \in \mathbb{R}^{n \times n}$ is the Laplacian matrix, $\mathbf{W} \in \mathbb{R}^{n \times n}$ is the affinity matrix and $\mathbf{D} = diag(\mathbf{W1}) \in \mathbb{R}^{n \times n}$ is the degree matrix.

**Compared Methods.** We compare MPEC-EPM and MPEC-ADM against 5 methods on the '4gauss' data set (see Figure 1). (i) LP relaxation simply relaxes the binary

constraint to $-1 \le \mathbf{x} \le 1$ and solves the following problem:

$$\mathbf{LP}: \quad \min_{\mathbf{x}} \mathbf{x}^T \mathbf{L} \mathbf{x}, \ s.t. \ \mathbf{x}^T \mathbf{1} = 0, \ -1 \le \mathbf{x} \le 1$$

(ii) Ratio Cut (RCUT) and Normalize Cut (NCUT) relax the binary constraint to $\|\mathbf{x}\|_2^2 = n$ and solve the following problems [51, 16]:

$$\mathbf{RCut}: \quad \min_{\mathbf{x}} \mathbf{x}^T \mathbf{L} \mathbf{x}, \ s.t. \ \langle \mathbf{x}, \mathbf{1} \rangle = 0, \ \|\mathbf{x}\|_2^2 = n$$

$$\mathbf{NCut}: \quad \min_{\mathbf{x}} \mathbf{x}^T \bar{\mathbf{L}} \mathbf{x}, \ s.t. \ \langle \mathbf{x}, \mathbf{D}^{1/2} \mathbf{1} \rangle = 0, \ \|\mathbf{x}\|_2^2 = n$$

where $\bar{\mathbf{L}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$. The optimal solution of RCut (or NCut) is the second smallest eigenvectors of $\mathbf{L}$ (or $\bar{\mathbf{L}}$), see e.g. [35, 60]. (iv) SDP relaxation solves the following convex optimization problem [15]:

$$\mathbf{SDP}: \quad \min_{\mathbf{X}} \langle \mathbf{L}, \mathbf{X} \rangle, \ s.t. \ diag(\mathbf{X}) = \mathbf{1}, \ \langle \mathbf{X}, \mathbf{11}^T \rangle = 0$$

---

[15]For SDP method, we use the randomized rounding strategy in [24, 60] to get a discrete solution from $\mathbf{X}$. Specifically, we sample a random vector $\mathbf{x} \in \mathbb{R}^n$ from a Gaussian distribution with mean 0 and covariance $\mathbf{X}$, and perform $\mathbf{x}^* = \text{sign}(\mathbf{x} - \text{median}(\mathbf{x}))$. This process is repeated many times and the final solution with the largest objective value is selected.

(a) BNCUT, $f = 89.09$    (b) BNCUT, $f = 23.08$    (c) BNCUT, $f = 26.67$    (d) BNCUT, $f = 56.02$    (e) BNCUT, $f = 75.09$

(f) LP, $f = 65.70$    (g) LP, $f = 27.53$    (h) LP, $f = 14.81$    (i) LP, $f = 27.50$    (j) LP, $f = 6.81$

(k) SDP, $f = 64.40$    (l) SDP, $f = 19.01$    (m) SDP, $f = 14.81$    (n) SDP, $f = 26.97$    (o) SDP, $f = 6.81$

(p) L2box-ADMM, $f = 56.65$    (q) L2box-ADMM, $f = 19.70$    (r) L2box-ADMM, $f = 14.32$    (s) L2box-ADMM, $f = 25.77$    (t) L2box-ADMM, $f = 6.81$

(u) MPEC-EPM, $f = 49.56$    (v) MPEC-EPM, $f = 16.41$    (w) MPEC-EPM, $f = 14.81$    (x) MPEC-EPM, $f = 25.77$    (y) MPEC-EPM, $f = 6.81$

(z) MPEC-ADM, $f = 50.46$    (aa) MPEC-ADM, $f = 16.41$    (ab) MPEC-ADM, $f = 13.50$    (ac) MPEC-ADM, $f = 24.97$    (ad) MPEC-ADM, $f = 6.20$
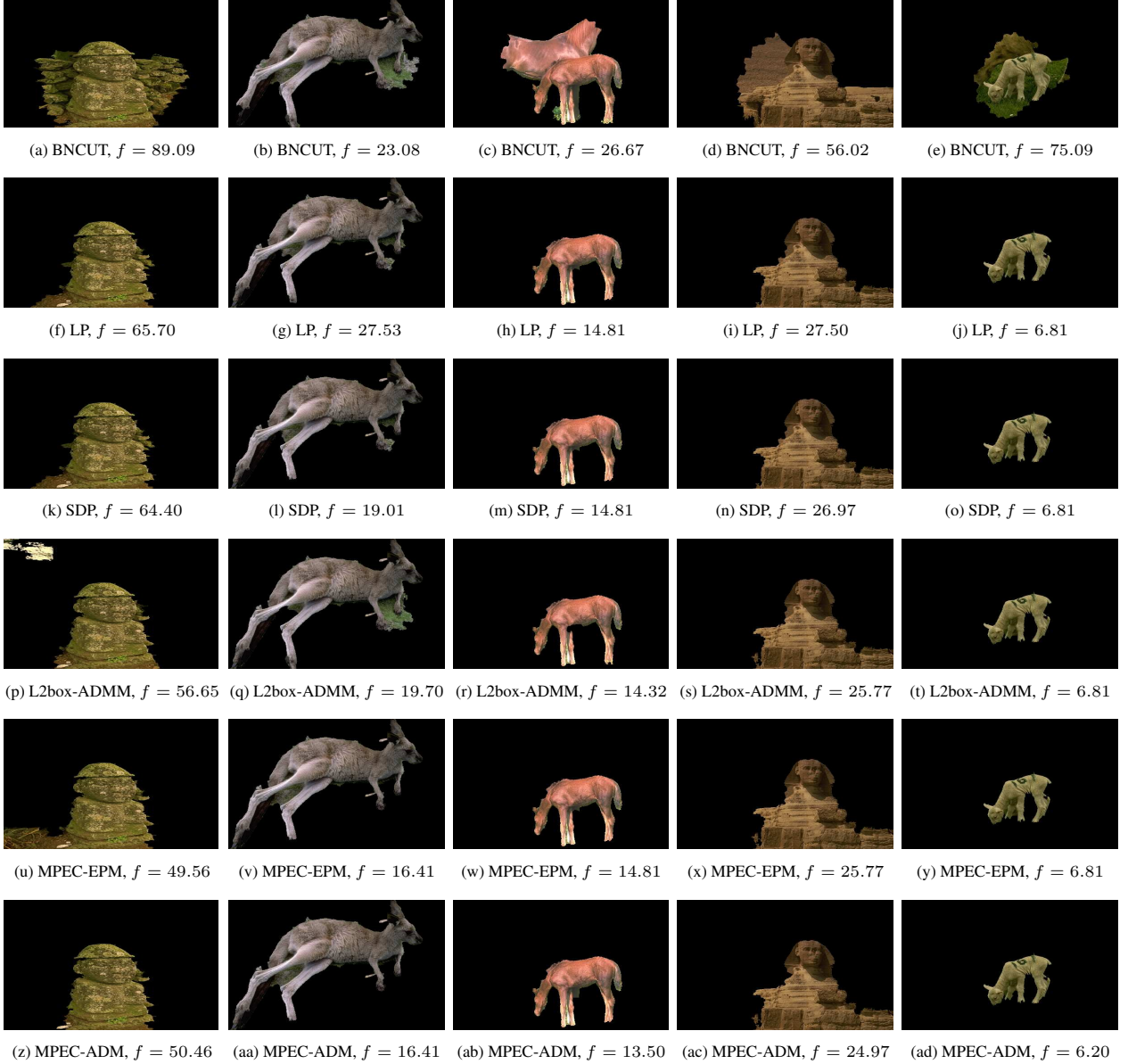
Figure 3: Images used in our constrained image segmentation experiments.

Finally, we also compare with L2box-ADMM [63] which applies ADMM directly to the $\ell_2$ box non-separable reformulation:

$$\min_{\mathbf{x}} \ \mathbf{x}^T \mathbf{L} \mathbf{x}, \ s.t. \ \mathbf{x}^T \mathbf{1} = 0, \ -\mathbf{1} \le \mathbf{x} \le \mathbf{1}, \ \|\mathbf{x}\|_2^2 = n.$$

We remark that it is a splitting method that introduces auxiliary variables to separate the two constrained set and then performs block coordinate descend on each variable.

**Experimental Results.** Several observations can be drawn from Figure 1. (i) The LP, RCUT and NCUT relaxation methods fail to appropriately separate the '4gauss' data set and they result in large objective values. (ii) SD-

P relaxation provides a good approximation and achieves a lower objective value than LP, RCUT, NCUT and L2box-ADMM. (iii) The proposed methods MPEC-EPM and MPEC-ADM achieve the same lowest objective values among all the compared methods.

## 5.2. Constrained Image Segmentation

In graph-based partition, image is modeled as a weighted undirected graph where nodes corresponds to pixels (or pixel regions) and edges encode feature similarities between the node pairs. Image segmentation can be treated as seeking a partition $\mathbf{x} \in \{-1, +1\}^n$ to cut the edges of the graph

with minimal weights. Prior label information on the vertices of the graph can be incorporated to improve performance, leading to the following optimization problem:

$$\min_{\mathbf{x}\in\{-1,+1\}^n} \mathbf{x}^T\mathbf{L}\mathbf{x}, \ s.t. \ \mathbf{x}_F = 1, \ \mathbf{x}_B = -1$$

where $F$ and $B$ denote the index of foreground pixels set and background pixels set, respectively; $\mathbf{L}\in\mathbb{R}^{n\times n}$ is the graph Laplacian matrix. Since SDP method can not solve large scale image segmentation problems, we over-segment the images into SLIC pixel regions using the 'VLFeat' toolbox [58]. The affinity matrix $\mathbf{W}$ is constructed based on the color similarities and spatial adjacencies between pixel regions.

**Compared Methods.** We compare MPEC-EPM and MPEC-ADM against 4 methods on the Weizman horses and MSRC datasets [60] (see Figure 2). (i) Biased normalized cut (BNCut) [16] extends Normalized Cut [51] to encode the labelled foreground pixels as a quadratic constraint on the solution $\mathbf{x}$. The solution of BNCut is a linear combination of the eigenvectors of normalized Laplacian matrix [51]. (ii) LP relaxation simply replaces the binary constraint with a soft constraint and solves a quadratic programming problem: $\min_{\mathbf{x}} \mathbf{x}^T\mathbf{L}\mathbf{x}, \ s.t. \ -\mathbf{1}\leq\mathbf{x}\leq\mathbf{1}, \ \mathbf{x}_F = 1, \ \mathbf{x}_B = -1$. (iii) SDP relaxation method considers the following optimization problem:

$$\min_{\mathbf{X}\succeq 0} \langle\mathbf{L},\mathbf{X}\rangle, \ s.t. \ diag(\mathbf{X}) = \mathbf{1}, \ \mathbf{X}_I = 1, \ \mathbf{X}_J = -1$$

with $\mathbf{X}\in\mathbb{R}^{n\times n}$, and $I$ and $J$ are the index pairs of similarity and dissimilarity, respectively. Therefore, it contains $n + \binom{2}{|B|} + \binom{2}{|F|} + |B|\cdot|F|$ linear equality constraints. We use 'cvx' optimization software [25] to solve this problem. (iv) L2box-ADMM considers the following optimization problem: $\min_{\mathbf{x}} \mathbf{x}^T\mathbf{L}\mathbf{x}, \ s.t. \ -\mathbf{1}\leq\mathbf{x}\leq\mathbf{1}, \ \mathbf{x}_F = 1, \ \mathbf{x}_B = -1, \ \|\mathbf{x}\|_2^2 = n$.

**Experimental Results.** Several observations can be drawn from Figure 3. (i) LP relaxation generates better image segmentation results than BNCUT except in the second image. Moreover, we found that a lower objective value does not always necessarily result in better view for image segmentation. We argue that the parameter in constructing the similarity graph is responsible for this result. (ii) SDP method and L2box-ADMM method generate better solutions than BNCUT and LP. (iii) The proposed MPEC-EPM and MPEC-ADM generally obtain lower objective values and outperform all the other compared methods.

### 5.3. Dense Subgraph Discovery

Dense subgraphs discovery [50, 20, 68] is a fundamental graph-theoretic problem, as it captures numerous graph mining applications, such as community finding, regulatory motifs detection, and real-time story identification. It aims at finding the maximum density subgraph on $k$ vertices,

Table 2: The statistics of the web graph data sets used in our dense subgraph discovery experiments.

| Graph | # Nodes | # Arcs | Avg. Degree |
|---|---|---|---|
| wordassociation | 10617 | 72172 | 6.80 |
| enron | 69244 | 276143 | 3.99 |
| uk-2007-05 | 100000 | 3050615 | 30.51 |
| cnr-2000 | 325557 | 3216152 | 9.88 |
| dblp-2010 | 326186 | 1615400 | 4.95 |
| in-2004 | 1382908 | 16917053 | 12.23 |
| amazon-2008 | 735323 | 5158388 | 7.02 |
| dblp-2011 | 986324 | 6707236 | 6.80 |

which can be formulated as the following binary program:

$$\max_{\mathbf{x}\in\{0,1\}^n} \mathbf{x}^T\mathbf{W}\mathbf{x}, \ s.t. \ \mathbf{x}^T\mathbf{1} = k \qquad (14)$$

where $\mathbf{W}\in\mathbb{R}^{n\times n}$ is the adjacency matrix of the graph. Although the objective function in Eq (14) may not be convex, one can append an additional term $\lambda\mathbf{x}^T\mathbf{x}$ to the objective with a sufficiently large $\lambda$ such that $\lambda\mathbf{I} - \mathbf{W} \succeq 0$. This is equivalent to adding a constant to the objective since $\lambda\mathbf{x}^T\mathbf{x} = \lambda k$ in the effective domain. Therefore, we have the following optimization problem which is equivalent to Eq (14):

$$\min_{\mathbf{x}\in\{0,1\}^n} \mathbf{x}^T(\lambda\mathbf{I} - \mathbf{W})\mathbf{x}, \ s.t. \ \mathbf{x}^T\mathbf{1} = k$$

In the experiments, $\lambda$ is set to the largest eigenvalue of $\mathbf{W}$.

**Compared Methods.** We compare MPEC-EPM and MPEC-ADM against 5 methods on 8 datasets [16] (see Table 2). (i) Feige's greedy algorithm (GEIGE) [20] is included in our comparisons. This method is known to achieve the best approximation ratio for general $k$. (ii) Ravi's greedy algorithm (RAVI) [50] starts from a heaviest edge and repeatedly adds a vertex to the current subgraph to maximize the weight of the resulting new subgraph. It has asymptotic performance guarantee of $\pi/2$ when the weights satisfy the triangle inequality. (iii) LP relaxation solves a capped simplex problem by standard quadratic programming technique: $\min_{\mathbf{x}} \mathbf{x}^T(\lambda\mathbf{I}-\mathbf{W})\mathbf{x}, \ s.t. \ \mathbf{0}\leq\mathbf{x}\leq\mathbf{1}, \ \mathbf{x}^T\mathbf{1} = k$. (iii-i) L2box-ADMM solves a spherical constraint optimization problem: $\min_{\mathbf{x}} \mathbf{x}^T(\lambda\mathbf{I} - \mathbf{W})\mathbf{x}, \ s.t. \ \mathbf{0}\leq\mathbf{x}\leq\mathbf{1}, \ \mathbf{x}^T\mathbf{1} = k, \ \|2\mathbf{x} - 1\|_2^2 = n$. (iv) Truncated Power Method (TPM) [68] considers an iterative procedure that combines power iteration and hard-thresholding truncation. It works by greedily decreasing the objective while maintaining the desired binary property for the intermediate solutions. We use the code [17] provided by the authors. As suggested in [68], the initial solution is set to the indicator vector of the vertices with the top $k$ weighted degrees of the graph in our experiments.

---

[16] http://law.di.unimi.it/datasets.php
[17] https://sites.google.com/site/xtyuan1980/publications
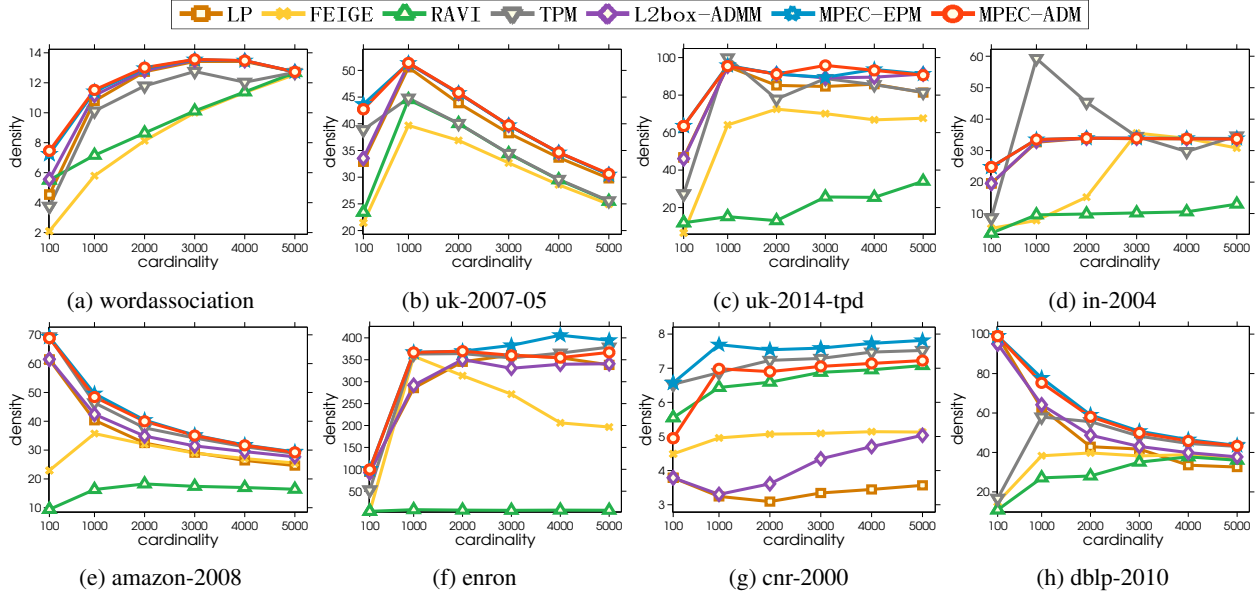
8

Figure 4: Experimental results for dense subgraph discovery.

**Experimental Results.** Several observations can be drawn from Figure 4. (i) Both FEIGE and RAVI generally fail to solve the dense subgraph discovery problem and they lead to solutions with low density. (ii) LP relaxation gives better performance than state-off-the-art technique TPM in some cases. (iii) L2box-ADMM outperforms LP relaxation for all cases but it generates unsatisfying accuracy in 'enron', 'cnr-2000' and 'dblp-2010'. (iv) Our proposed method MPEC-EPM and MPEC-ADM generally outperforms all the compared methods, while MPEC-EPM seems to present slightly better results than MPEC-ADM in this group of experiments.

## 5.4. Modularity Clustering

Modularity was first introduced in [45] as a performance measure for the quality of community structure found by a clustering algorithm. Given a modularity matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ with $\mathbf{Q}_{ij} = \mathbf{W}_{ij} - deg(\mathbf{v}_i)deg(\mathbf{v}_j)/(2m)$, modularity clustering can be formulated as the following optimization problem [44, 9, 14]:

$$\min_{\mathbf{X} \in \{-1, +1\}^{n \times n}} \frac{1}{8m} tr(\mathbf{X}^T \mathbf{Q} \mathbf{X}), \ s.t. \ \mathbf{X1} = (2 - k)\mathbf{1} \quad (15)$$

Observing that $\mathbf{Y} = (\mathbf{X} + \mathbf{1})/2 \in \mathbb{R}^{n \times n}$, we obtain $\mathbf{Y} \in \{0,1\}^{n \times n}$. Combining the linear constraint $\mathbf{X1} = (2 - k)\mathbf{1}$, we have $\mathbf{Y1} = \mathbf{1}$ and $\mathbf{XX}^T = (2\mathbf{Y} - \mathbf{1})(2\mathbf{Y} - \mathbf{1})^T = 4\mathbf{YY}^T - 3\mathbf{11}^T$. Based on these analyses, one can rewrite Eq (15) as the following equivalent optimization problem:

$$\min_{\mathbf{Y} \in \{0,1\}^{n \times n}} \frac{1}{8m} tr(\mathbf{Y}^T \mathbf{Q} \mathbf{Y}) + constant, \ s.t. \ \mathbf{Y1} = \mathbf{1}.$$

**Compared Methods.** We compare against 4 methods on 8 network data sets (See Table 3). (i) Iterative rounding algorithm (IRA) [14] in every stage solves a convex quadratic program and picks a fixed number of the vertices with largest values to assign to the cluster. However, such heuristic algorithm does not have any convergence guarantee. We use the code provided by the authors [18] and set the parameter $\rho = 0.5$ in this method. (ii) LP relaxation solves a capped simplex problem. (iii) L2box-ADMM solves a spherical constrained optimization problem. (iv) Iterative Hard Thresholding (IHT) considers setting the current solution to the indicator vector of its top-k entries while decreasing the objective function. Due to its suboptimal performance in our previous experiment, LP relaxation is used as its initialization.

**Experimental Results.** Several observations can be drawn from Figure 5. (i) IHT does not necessarily improve upon the LP relaxation method. (ii) IRA consistently outperforms IHT and LP in all the experiments. (iii) L2box-ADMM gives comparable result to IRA. (iv) Our proposed methods generally outperform IRA and L2box-ADM in the experiments.

## 5.5. Markov Random Fields

The Markov Random Field (MRF) optimization [8, 16, 31] is widely used in many labeling applications, including image restoration, edge detection, and image segmentation. Generally speaking, it involves solving the following prob-
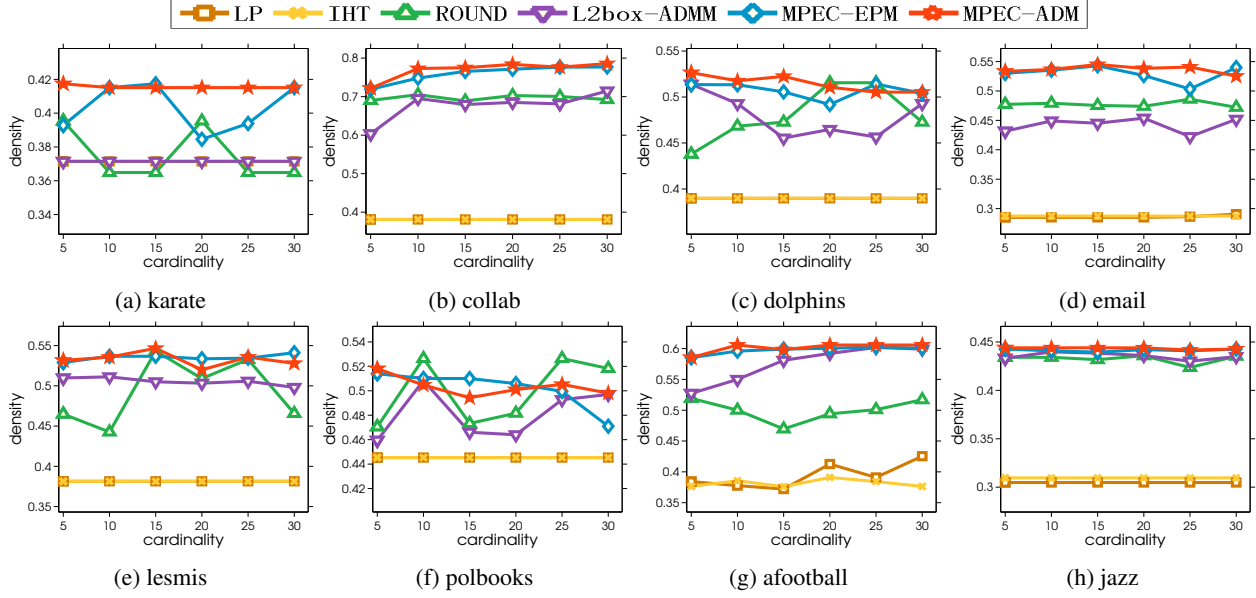
Figure 5: Experimental results for modularity clustering.

Table 3: The statistics of the web graph data sets used in our modularity clustering experiments.

| Graph | # Nodes | # Arcs | Avg. Degree |
|---|---|---|---|
| karate | 34 | 78 | 4.59 |
| collab | 235 | 415 | 3.53 |
| dolphins | 62 | 159 | 5.13 |
| email | 1133 | 5451 | 9.62 |
| lesmis | 77 | 820 | 21.30 |
| polbooks | 105 | 441 | 8.40 |
| afootball | 115 | 616 | 10.71 |
| jazz | 198 | 2742 | 27.70 |

Table 4: CPU time (in seconds) comparisons.

| Graph | LP | L2box-ADM | MPEC-EPM | MPEC-ADM |
|---|---|---|---|---|
| wordassociation | 1 | 7 | 2 | 13 |
| enron | 2 | 40 | 29 | 85 |
| uk-2007-05 | 6 | 75 | 65 | 77 |
| cnr-2000 | 16 | 210 | 209 | 245 |
| dblp-2010 | 15 | 234 | 282 | 253 |
| in-2004 | 79 | 834 | 1023 | 1301 |
| amazon-2008 | 49 | 501 | 586 | 846 |
| dblp-2011 | 59 | 554 | 621 | 1007 |

lem:

$$\min_{\mathbf{x}\in\{0,1\}^n} \frac{1}{2}\mathbf{x}^T\mathbf{L}\mathbf{x} + \mathbf{x}^T\mathbf{b}$$

where $\mathbf{b} \in \mathbb{R}^n$ is determined by the unary term defined for the graph and $\mathbf{L} \in \mathbb{R}^{n\times n}$ is the Laplacian matrix, which is based on the binary term relating pairs of graph nodes together. The quadratic term is usually considered a smoothness prior on the node labels.

**Compared Methods.** We perform image segmentation on the 'cat' image and 'flower' images. (i) Graph cut method [8] is included in our experiments. This method is known to achieve the global optimal solution for this specific class of binary problem. (ii) LP relaxation solves a box constrained quadratic programming problem. (iii) L2box-ADMM solves the $\ell_2$ box non-separable reformulation directly using classical alternating direction method of multipliers (ADMM) [63]. (iv) L0-QPM norm solves the semicontinuous $\ell_0$ norm reformulation of the binary optimization problem by quadratic penalty method [40, 67].

**Experimental Results.** Figure 6 demonstrates a qualitative result for image segmentation. MPEC-EPM and MPEC-ADM produce solutions that are very close to the globally optimal one. Moreover, both our methods achieve lower objectives than the other compared methods.

### 5.6. Convergence Curve and Computational Efficiency

This subsection demonstrates the convergence curve and computational efficiency of the proposed algorithms. We only report the results on dense subgraph discovery.

**Convergence Curve:** We demonstrate the convergence curve of the methods {LP, L2box-ADMM, MPEC-EPM, MPEC-ADM} for dense subgraph discovery on different data sets. As can be seen in Figure 7 and Figure 8, the proposed MPEC-based methods converges within 100 iterations. Moreover, we observe that the objective values generally decrease monotonically, and we attribute this to the greedy property of the penalty method for MPEC-EPM and monotone property of the dual variable $\rho$ update for
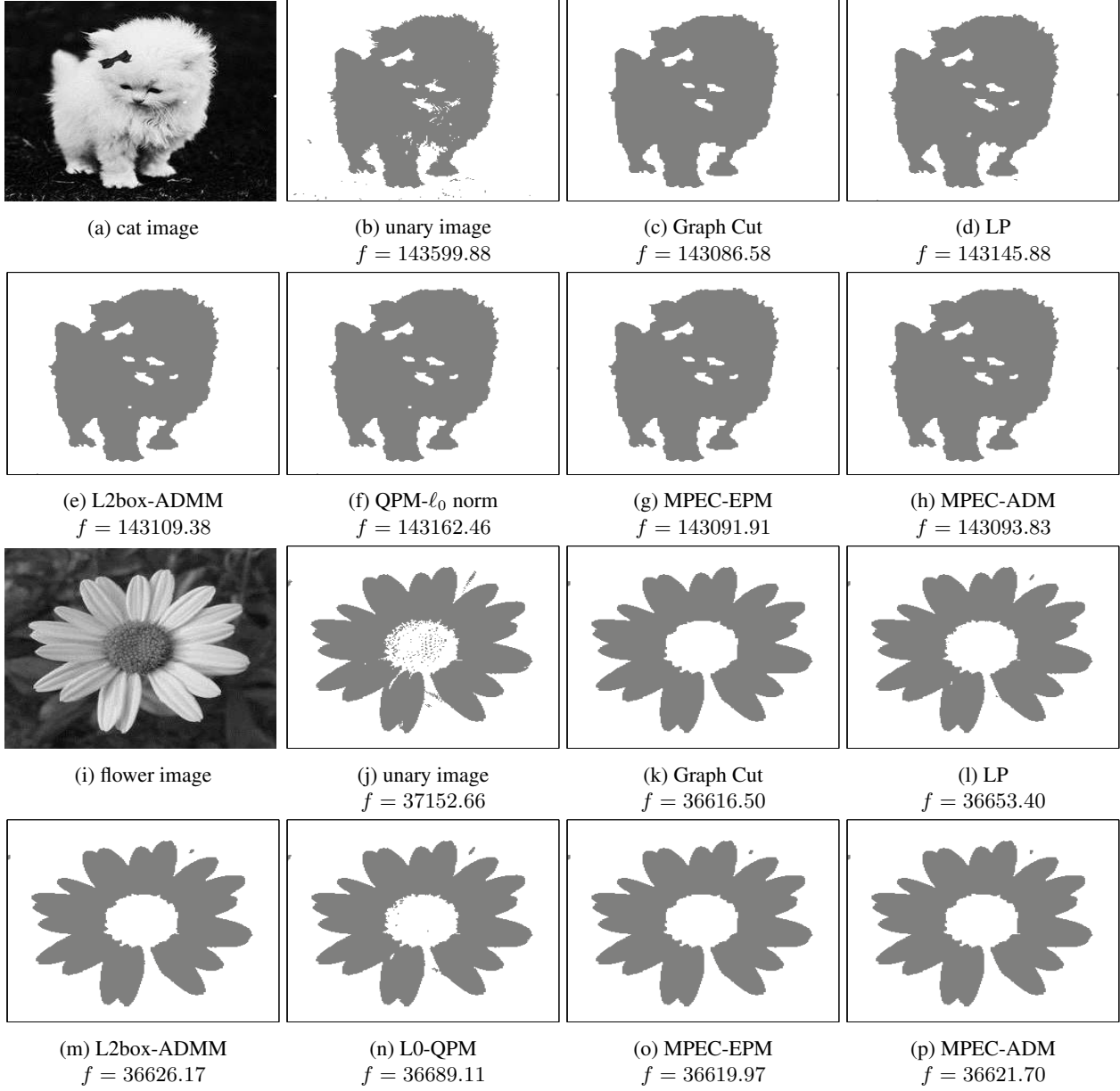
10

|                    |                    |                    |                    |
|--------------------|--------------------|--------------------|--------------------|
| (a) cat image      | (b) unary image $f = 143599.88$ | (c) Graph Cut $f = 143086.58$ | (d) LP $f = 143145.88$ |
| (e) L2box-ADMM $f = 143109.38$ | (f) QPM-$\ell_0$ norm $f = 143162.46$ | (g) MPEC-EPM $f = 143091.91$ | (h) MPEC-ADM $f = 143093.83$ |
| (i) flower image   | (j) unary image $f = 37152.66$ | (k) Graph Cut $f = 36616.50$ | (l) LP $f = 36653.40$ |
| (m) L2box-ADMM $f = 36626.17$ | (n) L0-QPM $f = 36689.11$ | (o) MPEC-EPM $f = 36619.97$ | (p) MPEC-ADM $f = 36621.70$ |

Figure 6: Markov random fields on 'cat' image and 'flower' images.

MPEC-ADM.

**Computational Efficiency:** We provide some running time comparisons for the methods {LP, L2box-ADMM, MPEC-EPM, MPEC-ADM} on different data sets with different $k \in \{100, 1000, 2000, 3000, 4000, 5000\}$. As can be seen in Table 4, even for the data set such as 'dblp-2011' that contains about one million nodes and 7 million edges, all the methods can terminate with in 15 minutes. Moreover, the runtime efficiency of our methods are several times slower than LP and comparable with and L2box-ADMM. This is expected, since (i) our methods MPEC-EPM and

MPEC-ADM need to call the LP procedure multiple times, and (ii) all the methods are all alternating methods and have the same computational complexity.

## 5.7. Some Implementation Details

This subsection presents some implementation details of MPEC-EPM and MPEC-ADM, including method of solving the **x**-subproblems and parameters setting of the algorithms.

The convex **x**-subproblems can be solved using Nesterov's projective gradient methods. For the application-
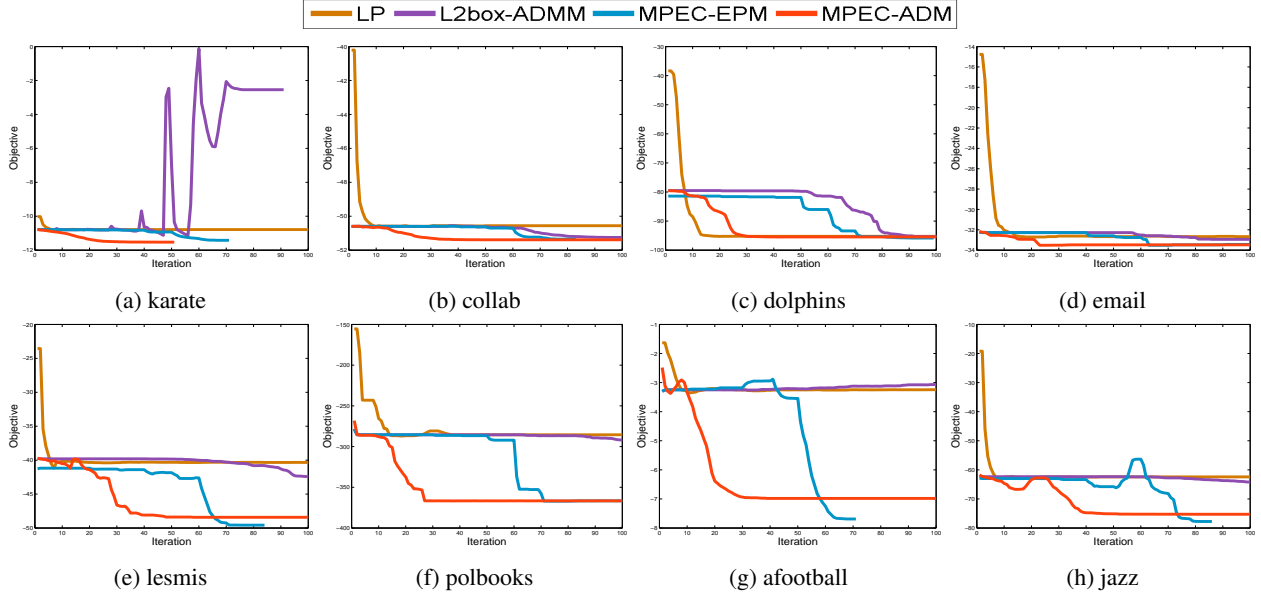
Figure 7: Convergence curve for dense subgraph discovery on different datasets with $k = 1000$.
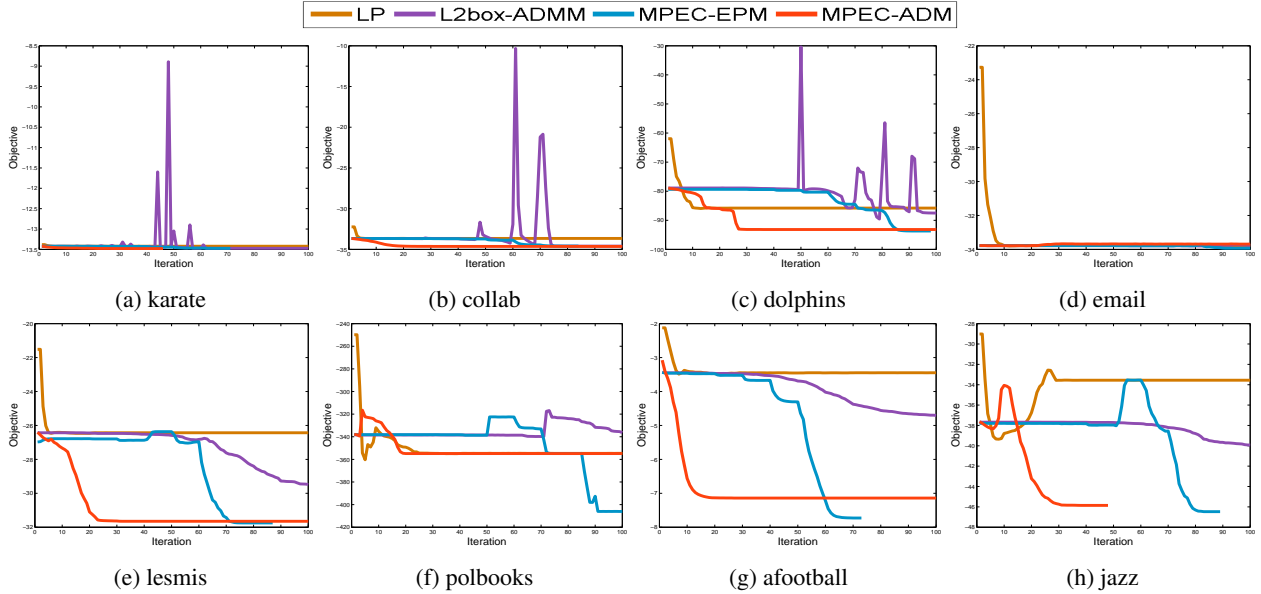


Figure 8: Convergence curve for dense subgraph discovery on different datasets with $k = 4000$.

s of constrained image segmentation and Markov random fields, the projection step involved in the **x**-subproblems is easy since it only contains box constraints [19]; for the applications of dense subgraph discovery, modularity clustering and graph bisection, the projection step can be hard since it contains an additional linear constraint besides box constraints (also known as capped simplex constraint) [20]. For-

tunately, this projection step can be solved by a break point search algorithm [28] exactly in $n \log(n)$ time. In our experiments, we use the Matlab implementation provided in the appendix of [67] .

The following parameters are used in our algorithms. For all methods, we use proximal gradient descent algorithm to solve the inner **x** subproblems. We stop the proximal gradient descent procedure when a relative change is smaller than $\epsilon = 10^{-5}$, i.e. $\|\mathbf{x}^{k+1} - \mathbf{x}^k\|/\|\mathbf{x}^k\| \leq \epsilon$, where $k$ is the iteration counter for the **x**-subproblem. In addition, we

---

[19]It solves: $\min_{\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}} \|\mathbf{a} - \mathbf{x}\|_2^2$, where $\mathbf{a}$ is given.

[20]It solves: $\min_{\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}, \ \mathbf{x}^T \mathbf{1} = k} \|\mathbf{a} - \mathbf{x}\|_2^2$, where $\mathbf{a}$ and $k$ are given.

set $\rho^0 = 0.01$, $T = 10$, $\sigma = \sqrt{10}$ for MPEC-EPM and set $\alpha^0 = 0.001$, $T = 10$, $\sigma = \sqrt{10}$ for MPEC-ADM. Finally, for L2box-ADM, we update the penalty parameter by a factor of $\sqrt{10}$ in every $T = 10$ iterations with initial value set to 0.1 [21].

## 6. Conclusions and Future Work

This paper presents a new class of continuous MPEC-based optimization methods to solve general binary optimization problems. Although the optimization problem is non-convex, we design two methods (exact penalty and alternating direction) to solve the equivalent problem. We also shed some theoretical lights to the equivalent formulations and optimization algorithms. Experimental results on binary problems demonstrate that our methods generally outperform existing solutions in terms of solution quality.

Our future work focuses on several directions. **(i)** We will investigate the optimality qualification of our multi-stage convex relaxation method for some specific objective functions, e.g., as is done in [24, 69, 13, 32]. **(ii)** We are also interested in extending the proposed algorithms to solve orthogonality and spherical optimization problems [62, 15] in computer vision and machine learning.

## Acknowledgments

## References

[1] Brendan P. Ames. Guaranteed clustering and biclustering via semidefinite programming. *Mathematical Programming*, 147(1-2):429–465, 2014. 1

[2] Brendan P. Ames. Guaranteed recovery of planted cliques and dense subgraphs by convex relaxation. *Journal of Optimization Theory and Applications*, 167(2):653–675, 2015. 1

[3] Brendan P. W. Ames and Stephen A. Vavasis. Nuclear norm minimization for the planted clique and biclique problems. *Mathematical Programming*, 129(1):69–89, 2011. 1

[4] Brendan P. W. Ames and Stephen A. Vavasis. Convex optimization for the planted k-disjoint-clique problem. *Mathematical Programming*, 143(1):299–337, 2014. 1

[5] Amir Beck and Yonina C. Eldar. Sparsity constrained nonlinear optimization: Optimality conditions and algorithms. *SIAM Journal on Optimization (SIOPT)*, 23(3):1480–1509, 2013. 2, 3

[6] Shujun Bi, Xiaolan Liu, and Shaohua Pan. Exact penalty decomposition method for zero-norm minimization based on mpec formulation. *SIAM Journal on Scientific Computing (SISC)*, 36(4):A1451–A1477, 2014. 3

[7] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2):459–494, 2014. 2, 4, 5

[8] Yuri Boykov, Olga Veksler, and Ramin Zabih. Fast approximate energy minimization via graph cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 23(11):1222–1239, 2001. 1, 9, 10

[9] Ulrik Brandes, Daniel Delling, Marco Gaertler, Robert Gorke, Martin Hoefer, Zoran Nikoloski, and Dorothea Wagner. On modularity clustering. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 20(2):172–188, 2008. 9

[10] Xavier Bresson, Xue-Cheng Tai, Tony F. Chan, and Arthur Szlam. Multi-class transductive learning based on $\ell_1$ relaxations of cheeger cut and mumford-shah-potts model. *Journal of Mathematical Imaging and Vision*, 49(1):191–201, 2014. 3

[11] Samuel Burer. On the copositive representation of binary and continuous nonconvex quadratic programs. *Mathematical Programming*, 120(2):479–495, 2009. 2, 3

[12] Samuel Burer. Optimizing a polyhedral-semidefinite relaxation of completely positive programs. *Mathematical Programming Computation*, 2(1):1–19, 2010. 2, 3

[13] Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015. 2, 13

[14] Emprise Y. K. Chan and Dit-Yan Yeung. A convex formulation of modularity maximization for community detection. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2218–2225, 2011. 1, 9

[15] Weiqiang Chen, Hui Ji, and Yanfei You. An augmented lagrangian method for l1-regularized optimization problems with orthogonality constraints. *SIAM*

---

[21] We tune this parameter in the range $\{0.01, 0.1, 1\}$ and find the value 0.1 generally gives comparable results.

*Journal on Scientific Computing (SISC)*, 38(4):B570–B592, 2016. 5, 13

[16] Timothee Cour and Jianbo Shi. Solving markov random fields with spectral relaxation. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 2, page 15, 2007. 2, 6, 8, 9

[17] Timothee Cour, Praveen Srinivasan, and Jianbo Shi. Balanced graph matching. *Neural Information Processing Systems (NIPS)*, 19:313, 2007. 1

[18] Mark A Davenport, Yaniv Plan, Ewout van den Berg, and Mary Wootters. 1-bit matrix completion. *Information and Inference*, 3(3):189–223, 2014. 1

[19] Marianna De Santis and Francesco Rinaldi. Continuous reformulations for zero–one programming problems. *Journal of Optimization Theory and Applications*, 153(1):75–84, 2012. 3

[20] Uriel Feige, David Peleg, and Guy Kortsarz. The dense k-subgraph problem. *Algorithmica*, 29(3):410–421, 2001. 8

[21] Fajwel Fogel, Rodolphe Jenatton, Francis R. Bach, and Alexandre d'Aspremont. Convex relaxations for permutation problems. *SIAM Journal on Matrix Analysis and Applications (SIMAX)*, 36(4):1465–1488, 2015. 1

[22] Walter Gander, Gene H Golub, and Urs von Matt. A constrained eigenvalue problem. *Linear Algebra and its applications*, 114:815–839, 1989. 5

[23] Bissan Ghaddar, Juan C. Vera, and Miguel F. Anjos. Second-order cone relaxations for binary quadratic polynomial programs. *SIAM Journal on Optimization (SIOPT)*, 21(1):391–414, 2011. 2

[24] Michel X. Goemans and David P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM*, 42(6):1115–1145, 1995. 1, 3, 6, 13

[25] Michael Grant and Stephen Boyd. CVX: Matlab software for disciplined convex programming, version 2.1. http://cvxr.com/cvx, March 2014. 3, 8

[26] Bingsheng He and Xiaoming Yuan. On the $\mathcal{O}(1/n)$ convergence rate of the douglas-rachford alternating direction method. *SIAM Journal on Numerical Analysis (SINUM)*, 50(2):700–709, 2012. 4, 18

[27] Lifang He, Chun-Ta Lu, Jiaqi Ma, Jianping Cao, Linlin Shen, and Philip S. Yu. Joint community and structural hole spanner detection via harmonic modularity. In *ACM SIGKDD Conferences on Knowledge Discovery and Data Mining (SIGKDD)*. ACM, 2016. 1

[28] R Helgason, J Kennington, and H Lall. A polynomially bounded algorithm for a singly constrained quadratic program. *Mathematical Programming*, 18(1):338–343, 1980. 12

[29] Cho-Jui Hsieh, Nagarajan Natarajan, and Inderjit S Dhillon. Pu learning for matrix completion. In *International Conference on Machine Learning (ICML)*, pages 2445–2453, 2015. 1, 2

[30] XM Hu and Daniel Ralph. Convergence of a penalty method for mathematical programming with complementarity constraints. *Journal of Optimization Theory and Applications*, 123(2):365–390, 2004. 4

[31] Qi-Xing Huang, Yuxin Chen, and Leonidas J. Guibas. Scalable semidefinite relaxation for maximum A posterior estimation. In *International Conference on Machine Learning (ICML)*, pages 64–72, 2014. 2, 3, 9

[32] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *ACM Symposium on Theory of Computing (STOC)*, pages 665–674, 2013. 2, 13

[33] Armand Joulin, Francis R. Bach, and Jean Ponce. Discriminative clustering for image co-segmentation. In *IEEE Conference on Computer Vision and Pattern (CVPR)*, pages 1943–1950, 2010. 1

[34] Bahman Kalantari and J. B. Rosen. Penalty for zero–one integer equivalent problem. *Mathematical Programming*, 24(1):229–232, 1982. 3

[35] Jens Keuchel, Christoph Schnorr, Christian Schellewald, and Daniel Cremers. Binary partitioning, perceptual grouping, and restoration with semidefinite programming. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 25(11):1364–1379, 2003. 1, 2, 6

[36] Nikos Komodakis and Georgios Tziritas. Approximate labeling via graph cuts based on linear programming. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 29(8):1436–1453, 2007. 2, 3, 4

[37] M. Pawan Kumar, Vladimir Kolmogorov, and Philip H. S. Torr. An analysis of convex relaxations for MAP estimation of discrete mrfs. *Journal of Machine Learning Research (JMLR)*, 10:71–106, 2009. 2, 3

[38] Guosheng Lin, Chunhua Shen, David Suter, and Anton van den Hengel. A general two-step approach to

learning-based hashing. In *International Conference on Computer Vision (ICCV)*, pages 2552–2559, 2013. 2

[39] Wei Liu, Jun Wang, Sanjiv Kumar, and Shih-Fu Chang. Hashing with graphs. In *International Conference on Machine Learning (ICML)*, pages 1–8, 2011. 1

[40] Zhaosong Lu and Yong Zhang. Sparse approximation via penalty decomposition methods. *SIAM Journal on Optimization (SIOPT)*, 23(4):2448–2478, 2013. 2, 3, 4, 10

[41] Zhi-Quan Luo, Jong-Shi Pang, and Daniel Ralph. *Mathematical programs with equilibrium constraints*. Cambridge University Press, 1996. 3, 4

[42] Walter Murray and Kien-Ming Ng. An algorithm for nonlinear optimization problems with binary variables. *Computational Optimization and Applications*, 47(2):257–288, 2010. 2, 3

[43] Y. E. Nesterov. *Introductory lectures on convex optimization: a basic course*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, 2003. 4

[44] Mark EJ Newman. Modularity and community structure in networks. *Proceedings of the national academy of sciences*, 103(23):8577–8582, 2006. 9

[45] Mark EJ Newman and Michelle Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2):026113, 2004. 9

[46] Carl Olsson, Anders P Eriksson, and Fredrik Kahl. Solving large scale binary quadratic problems: Spectral methods vs. semidefinite programming. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8. IEEE, 2007. 2

[47] Panos M. Pardalos and J. Ben Rosen. *Constrained Global Optimization: Algorithms and Applications*. Springer-Verlag New York, Inc., New York, NY, USA, 1987. 3

[48] M. Raghavachari. On connections between zero-one integer programming and concave programming under linear constraints. *Operations Research*, 17(4):680–684, 1969. 2, 3

[49] Daniel Ralph and Stephen J Wright. Some properties of regularization and penalization schemes for mpecs. *Optimization Methods and Software*, 19(5):527–556, 2004. 3, 4

[50] Sekharipuram S Ravi, Daniel J Rosenkrantz, and Giri K Tayi. Heuristic and special case algorithms for dispersion problems. *Operations Research*, 42(2):299–310, 1994. 8

[51] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 22(8):888–905, 2000. 1, 2, 6, 8

[52] Xinchu Shi, Haibin Ling, Junliang Xing, and Weiming Hu. Multi-target tracking by rank-1 tensor approximation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2387–2394, 2013. 1

[53] Fatemeh Shokrollahi Yancheshmeh, Ke Chen, and Joni-Kristian Kamarainen. Unsupervised visual alignment with similarity graphs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2901–2908, 2015. 1

[54] Ashish Shrivastava, Mohammad Rastegari, Sumit Shekhar, Rama Chellappa, and Larry S. Davis. Class consistent multi-modal fusion with binary features. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2282–2291, 2015. 1

[55] Sonja Steffensen and Michael Ulbrich. A new relaxation scheme for mathematical programs with equilibrium constraints. *SIAM Journal on Optimization (SIOPT)*, 20(5):2504–2539, 2010. 3, 4

[56] Alexander Toshev, Jianbo Shi, and Kostas Daniilidis. Image matching via saliency region correspondences. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8. IEEE, 2007. 1

[57] Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109(3):475–494, 2001. 2, 4, 5

[58] Andrea Vedaldi and Brian Fulkerson. Vlfeat: An open and portable library of computer vision algorithms. In *ACM International Conference on Multimedia (ACM MM)*, pages 1469–1472, 2010. 8

[59] Jun Wang, Wei Liu, Sanjiv Kumar, and Shih-Fu Chang. Learning to hash for indexing big data - a survey. *Proceedings of the IEEE*, 104(1):34–57, 2016. 1

[60] P. Wang, C. Shen, A. van den Hengel, and P. Torr. Large-scale binary quadratic optimization using semidefinite relaxation and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2016. 1, 2, 6, 8

15

[61] Zaiwen Wen, Donald Goldfarb, and Wotao Yin. Alternating direction augmented lagrangian methods for semidefinite programming. *Mathematical Programming Computation*, 2(3):203–230, 2010. 2, 4

[62] Zaiwen Wen and Wotao Yin. A feasible method for optimization with orthogonality constraints. *Mathematical Programming*, 142(1-2):397–434, 2013. 5, 13

[63] Baoyuan Wu and Bernard Ghanem. $\ell_p$-box ADMM: A versatile framework for integer programming. In *arXiv preprint*, 2016. 2, 3, 7, 10

[64] Junchi Yan, Chao Zhang, Hongyuan Zha, Wei Liu, Xiaokang Yang, and Stephen M Chu. Discrete hypergraph matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1520–1528, 2015. 1

[65] Ganzhao Yuan and Bernard Ghanem. $\ell_0 tv$: A new method for image restoration in the presence of impulse noise. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5369–5377, 2015. 3, 4

[66] Ganzhao Yuan and Bernard Ghanem. A proximal alternating direction method for semi-definite rank minimization. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2016. 3, 4

[67] Ganzhao Yuan and Bernard Ghanem. Sparsity constrained minimization via mathematical programming with equilibrium constraints. In *arXiv preprint*, 2016. 2, 3, 4, 10, 12

[68] Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research (JMLR)*, 14(1):899–925, 2013. 1, 2, 3, 8

[69] Tong Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *The Journal of Machine Learning Research (JMLR)*, 11:1081–1107, 2010. 2, 3, 13

[70] Zhongyuan Zhang, Tao Li, Chris Ding, and Xiangsun Zhang. Binary matrix factorization with applications. In *IEEE International Conference on Data Mining (ICDM)*, pages 391–400. IEEE, 2007. 2, 3

[71] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. 5

# Appendix

## A. Proofs of Equivalent Reformulations

This section presents the proofs of the equivalent separable/non-separable reformulations for binary constraint which are claimed in Table 1. Note that separable MPEC has $n$ complementarity constraints and non-separable MPEC has one complementarity constraint. The terms separable and non-separable are related to whether the constraints can be decomposed to independent components.

**Lemma 1.** $\ell_2$ *box non-separable MPEC. We define*

$$\Theta \triangleq \{(\mathbf{x}, \mathbf{v}) \mid \mathbf{x}^T \mathbf{v} = n, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_2^2 \leq n\}. \quad (16)$$

*Assume that* $(\mathbf{x}, \mathbf{v}) \in \Theta$, *then we have* $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$, *and* $\mathbf{x} = \mathbf{v}$.

*Proof.* (i) Firstly, we prove that $\mathbf{x} \in \{-1, +1\}^n$. Using the definition of $\Theta$ and the Cauchy-Schwarz Inequality, we have: $n = \mathbf{x}^T\mathbf{v} \leq \|\mathbf{x}\|_2\|\mathbf{v}\|_2 \leq \|\mathbf{x}\|_2\sqrt{n} = \sqrt{n\mathbf{x}^T\mathbf{x}} \leq \sqrt{n\|\mathbf{x}\|_1\|\mathbf{x}\|_\infty} \leq \sqrt{n\|\mathbf{x}\|_1}$. Thus, we obtain $\|\mathbf{x}\|_1 \geq n$. We define $\mathbf{z} = |\mathbf{x}|$. Combining $\|\mathbf{x}\|_\infty \leq 1$, we have the following constraint sets for $\mathbf{z}$: $\sum_i \mathbf{z}_i \geq n$, $\mathbf{0} \leq \mathbf{z} \leq \mathbf{1}$. Therefore, we have $\mathbf{z} = \mathbf{1}$ and it holds that $\mathbf{x} \in \{-1, +1\}^n$. (ii) Secondly, we prove that $\mathbf{v} \in \{-1, +1\}^n$. We have:

$$n = \mathbf{x}^T\mathbf{v} \leq \|\mathbf{x}\|_\infty\|\mathbf{v}\|_1 \leq \|\mathbf{v}\|_1 = |\mathbf{v}|^T\mathbf{1} \leq \|\mathbf{v}\|_2\|\mathbf{1}\|_2 \ (17)$$

Thus, we obtain $\|\mathbf{v}\|_2 \geq \sqrt{n}$. Combining $\|\mathbf{v}\|_2^2 \leq n$, we have $\|\mathbf{v}\|_2 = \sqrt{n}$ and $\|\mathbf{v}\|_2\|\mathbf{1}\|_2 = n$. By the Squeeze Theorem, all the equalities in Eq (17) hold automatically. Using the equality condition for Cauchy-Schwarz Inequality, we have $|\mathbf{v}| = \mathbf{1}$ and it holds that $\mathbf{v} \in \{-1, +1\}^n$.

(iii) Finally, since $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$, and $\langle \mathbf{x}, \mathbf{v} \rangle = n$, we obtain $\mathbf{x} = \mathbf{v}$. $\square$

**Lemma 2.** $\ell_\infty$ *box non-separable MPEC. We define*

$$\Pi \triangleq \{(\mathbf{x}, \mathbf{v}) \mid \mathbf{x}^T \mathbf{v} = n, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_\infty \leq 1\}.$$

*Assume that* $(\mathbf{x}, \mathbf{v}) \in \Pi$, *then we have* $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$, *and* $\mathbf{x} = \mathbf{v}$.

*Proof.* The proof of this lemma is similar to that of Lemma 1. Firstly, we prove that $\mathbf{x} \in \{-1, +1\}^n$. We have: $n = \mathbf{x}^T\mathbf{v} \leq \|\mathbf{x}\|_1 \cdot \|\mathbf{v}\|_\infty \leq \|\mathbf{x}\|_1$. Thus, we obtain $\|\mathbf{x}\|_1 \geq n$. We define $\mathbf{z} = |\mathbf{x}|$. Combining $\|\mathbf{x}\|_\infty \leq 1$, we have the following constraint sets for $\mathbf{z}$: $\sum_i \mathbf{z}_i \geq n$, $\mathbf{0} \leq \mathbf{z} \leq \mathbf{1}$. Therefore, we have $\mathbf{z} = \mathbf{1}$ and it holds that $\mathbf{x} \in \{-1, +1\}^n$. Secondly, using the same methodology, we can prove that $\mathbf{v} \in \{-1, +1\}^n$. Finally, we have $\mathbf{x} = \mathbf{v}$ since $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$ and $\langle \mathbf{x}, \mathbf{v} \rangle = n$. $\square$

**Lemma 3.** $\ell_\infty$ *box separable MPEC. We define*

$$\Psi \triangleq \{(\mathbf{x}, \mathbf{v}) \mid \mathbf{x} \odot \mathbf{v} = \mathbf{1}, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_\infty \leq 1\}.$$

*Assume that* $(\mathbf{x}, \mathbf{v}) \in \Psi$, *then we have* $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$, *and* $\mathbf{x} = \mathbf{v}$.

*Proof.* We observe that all the constraints in $\Psi$ can be decomposed into $n$ independent components. We now focus on $i$th component. (i) Assuming that $\mathbf{x}_i \geq 0$, we have $0 \leq \mathbf{x}_i \leq 1$ and $0 \leq 1/\mathbf{x}_i \leq 1$, we have $\mathbf{x}_i = 1$ and $\mathbf{v}_i = 1/\mathbf{x}_i = 1$. (ii) Assuming that $\mathbf{x}_i \leq 0$, we have $-1 \leq \mathbf{x}_i \leq 0$ and $-1 \leq 1/\mathbf{x}_i \leq 0$, we have $\mathbf{x}_i = -1$ and $\mathbf{v}_i = 1/\mathbf{x}_i = -1$. This finishes the proof. $\square$

**Lemma 4.** $\ell_2$ *box separable MPEC. We define*

$$\Upsilon \triangleq \{(\mathbf{x}, \mathbf{v}) \mid \mathbf{x} \odot \mathbf{v} = 1, \ \|\mathbf{x}\|_\infty \leq 1, \ \|\mathbf{v}\|_2^2 \leq n\}$$

*Assume that* $(\mathbf{x}, \mathbf{v}) \in \Upsilon$, *then we have* $\mathbf{x} \in \{-1, +1\}^n$, $\mathbf{v} \in \{-1, +1\}^n$, *and* $\mathbf{x} = \mathbf{v}$.

*Proof.* We notice that $\mathbf{v} = \frac{1}{\mathbf{x}}$. We define $\mathbf{z} = |\mathbf{v}|$. Combining $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$, we obtain $\mathbf{z} \geq \mathbf{1}$. We have the following constraint sets for $\mathbf{z}$: $\mathbf{z} \geq \mathbf{1}$, $\mathbf{z}^T\mathbf{z} \leq n$. Therefore, we have $\mathbf{z} = \mathbf{1}$. Finally, we achieve $\mathbf{v} \in \{-1, +1\}^n$ and $\mathbf{x} = \frac{1}{\mathbf{v}} = \mathbf{v}$. $\square$

**Lemma 5.** $\ell_2$ *box non-separable reformulation. The following equivalence holds:*

$$\{-1 + 1\}^n \Leftrightarrow \{\mathbf{x} \mid -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ \|\mathbf{x}\|_2^2 = n\}$$

*Proof.* First, it holds that: $n = \mathbf{x}^T\mathbf{x} \leq \|\mathbf{x}\|_1 \|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1$. Therefore, we have $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$, $\|\mathbf{x}\|_1 \geq n$. Note that these constraint sets are symmetric. Letting $\mathbf{z} = |\mathbf{x}|$, we obtain: $\mathbf{0} \leq \mathbf{z} \leq \mathbf{1}$, $\sum_i \mathbf{z}_i \geq n$. Thus, we have $\mathbf{z} = \mathbf{1}$ and it holds that $\mathbf{x} \in \{-1, +1\}^n$. $\square$

## B. Proofs of Convergence Theorems

This section provides the proofs of convergence theorems of Algorithm 1 and Algorithm 2.

The following lemma is very useful in our proofs.

**Lemma 6.** *Let* $\mathbf{x} \in \mathbb{R}^n$ *be an arbitrary vector with* $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. *We assume* $\text{sign}(\mathbf{x}) \neq \mathbf{x}$ *and define* $\text{sign}(0) = \pm 1$. *The following inequalities hold:*

$$h(\mathbf{x}) \triangleq \frac{n - \sqrt{n}\|\mathbf{x}\|_2}{\|\text{sign}(\mathbf{x}) - \mathbf{x}\|_2} > n - \sqrt{n^2 - n} > 1/2 \quad (18)$$

*Proof.* (i) We prove the first inequality in Eq (18). We define $\mathcal{N}(\mathbf{x})$ as the number of $\pm 1$ binary variables in $\mathbf{x}$, i.e., $\mathcal{N}(\mathbf{x}) \triangleq \#(|\mathbf{x}| = 1)$. Clearly, the objective function $h(\mathbf{x})$ decreases as $\mathcal{N}(\mathbf{x})$ increases. Note that $\mathcal{N}(\mathbf{x}) \neq n$, since

otherwise it violates the assumption that $\text{sign}(\mathbf{x}) \neq \mathbf{x}$. We consider the objective value $h(\mathbf{x})$ when $\mathcal{N}(\mathbf{x}) = n - 1$. In this situation, there exists only one coordinate such that $\text{sign}(\mathbf{x}_i) \neq \mathbf{x}_i$ with $\mathbf{x}_i = \pm\delta$, $0 < \delta < 1$ and the remaining coordinates take binary variable in $\{-1, +1\}$. Note that $\delta \neq 0$ and $\delta \neq 1$, since otherwise it also violates the assumption that $\text{sign}(\mathbf{x}) \neq \mathbf{x}$. Therefore, we derive the following inequalities:

$$
\begin{aligned}
& \frac{n - \sqrt{n}\|\mathbf{x}\|_2}{\|\text{sign}(\mathbf{x}) - \mathbf{x}\|_2} \\
> \ & \frac{n - \sqrt{n}\sqrt{(n-1) + \delta^2}}{\sqrt{(1-\delta)^2}} \\
\geq \ & \frac{n - \sqrt{n}(\sqrt{n-1} + \delta)}{(1-\delta)} \\
= \ & \frac{n - \sqrt{n}\sqrt{n-1}}{(1-\delta)} + \frac{\sqrt{n}\delta}{(1-\delta)} \\
> \ & \frac{n - \sqrt{n}\sqrt{n-1}}{1} + 0
\end{aligned}
\tag{19}
$$

where the second step uses $\sqrt{a + b} \leq \sqrt{a} + \sqrt{b}$, $\forall a, b > 0$; the last step uses the fact that $0 < \delta < 1$. Since the lower bound in Eq (19) can be applied to an arbitrary vector, we finish the proof of the first inequality.

(ii) We prove the second inequality in Eq (18). We have the following results:

$$
\begin{aligned}
& 1/4 > 0 \\
\Rightarrow \ & n^2 - n + 1/4 > n^2 - n \\
\Rightarrow \ & (n - 1/2)^2 > n^2 - n \\
\Rightarrow \ & (n - 1/2) > \sqrt{n^2 - n} \\
\Rightarrow \ & n - \sqrt{n^2 - n} > 1/2
\end{aligned}
$$

Hereby we finish the proof of this lemma. $\square$

The following lemma is useful in establishing the exactness property of the penalty function in Algorithm 1.

**Lemma 7.** *Consider the following optimization problem:*

$$(\mathbf{x}_\rho^*, \mathbf{v}_\rho^*) = \arg \min_{-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \ \|\mathbf{v}\|_2^2 \leq n, \ \mathbf{x} \in \Omega} \mathcal{J}_\rho(\mathbf{x}, \mathbf{v}). \quad (20)$$

*Assume that* $f(\cdot)$ *is a* $L$-*Lipschitz continuous convex function on* $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. *When* $\rho > 2L$, $\langle \mathbf{x}_\rho^*, \mathbf{v}_\rho^* \rangle = n$ *will be achieved for any local optimal solution of Eq (20).*

*Proof.* First of all, we focus on the $\mathbf{v}$-subproblem in the optimization problem in Eq (20).

$$\mathbf{v}_\rho^* = \arg \min_{\mathbf{v}} \ -\mathbf{x}^T\mathbf{v}, \ s.t. \ \|\mathbf{v}\|_2^2 \leq n \quad (21)$$

17

Assume that $\mathbf{x}_\rho^* \neq \mathbf{0}$, we have $\mathbf{v}_\rho^* = \sqrt{n} \cdot \mathbf{x}_\rho^*/\|\mathbf{x}_\rho^*\|_2$. Then the biconvex optimization problem reduces to the following:

$$\mathbf{x}_\rho^* = \arg \min_{-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}} p(\mathbf{x}) \triangleq f(\mathbf{x}) + \rho(n - \sqrt{n}\|\mathbf{x}\|_2) \quad (22)$$

For any $\mathbf{x}_\rho^* \in \Omega$, we derive the following inequalities:

$$
\begin{aligned}
& 0.5\rho\|\text{sign}(\mathbf{x}_\rho^*) - \mathbf{x}_\rho^*\|_2 \\
\leq\ & \rho(n - \sqrt{n}\|\mathbf{x}_\rho^*\|_2) \\
=\ & [\rho(n - \sqrt{n}\|\mathbf{x}_\rho^*\|_2) + f(\mathbf{x}_\rho^*)] - f(\mathbf{x}_\rho^*) \\
\leq\ & [\rho(n - \sqrt{n}\|\text{sign}(\mathbf{x}_\rho^*)\|_2) + \mathrm{f}(\text{sign}(\mathbf{x}_\rho^*))] - \mathrm{f}(\mathbf{x}_\rho^*) \\
=\ & f(\text{sign}(\mathbf{x}_\rho^*)) - \mathrm{f}(\mathbf{x}_\rho^*) \\
=\ & L\|\text{sign}(\mathbf{x}_\rho^*) - \mathbf{x}_\rho^*\|_2 \quad (23)
\end{aligned}
$$

where the first step uses Lemma 6 that $\|\text{sign}(\mathbf{x}) - \mathbf{x}\|_2 \leq 2(\mathrm{n} - \sqrt{n}\|\mathbf{x}\|_2)$ for any $\mathbf{x}$ in $\|\mathbf{x}\|_\infty \leq 1$; the third steps uses the optimality of $\mathbf{x}_\rho^*$ in Eq (22) that $p(\mathbf{x}_\rho^*) \leq p(\mathbf{y})$ for any $\mathbf{y}$ with $-\mathbf{1} \leq \mathbf{y} \leq \mathbf{1}$, $\mathbf{y} \in \Omega$; the fourth step uses the fact that $\text{sign}(\mathbf{x}_\rho) \in \{-1, +1\}^\mathrm{n}$ and $\sqrt{n}\|\text{sign}(\mathbf{x}_\rho)\|_2 = \mathrm{n}$; the last step uses the Lipschitz continuity of $f(\cdot)$.

From Eq (23), we have $\|\mathbf{x}_\rho^* - \text{sign}(\mathbf{x}_\rho^*)\|_2 \cdot (\rho - 2\mathrm{L}) \leq 0$. Since $\rho - 2L > 0$, we conclude that it always holds that $\|\mathbf{x}_\rho^* - \text{sign}(\mathbf{x}_\rho^*)\|_2 = 0$. Thus, $\mathbf{x}_\rho^* \in \{-1, +1\}^n$. Finally, we have $\mathbf{x}_\rho^* = \sqrt{n} \cdot \mathbf{x}_\rho^*/\|\mathbf{x}_\rho^*\|_2 = \mathbf{v}_\rho^*$ and $\langle \mathbf{x}_\rho^*, \mathbf{v}_\rho^* \rangle = n$. $\qquad \square$

**Theorem 1.** *Exactness of the Penalty Function. Assume that $f(\cdot)$ is a L-Lipschitz continuous convex function on $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. When $\rho > 2L$, the biconvex optimization $\min_{\mathbf{x}, \mathbf{v}} \mathcal{J}_\rho(\mathbf{x}, \mathbf{v})$, $s.t. -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$, $\|\mathbf{v}\|_2^2 \leq n$ in Eq (3) has the same local and global minima with the original problem in Eq (2).*

*Proof.* We let $\mathbf{x}^*$ be any global minimizer of Eq (2) and $(\mathbf{x}_\rho^*, \mathbf{v}_\rho^*)$ be any global minimizer of Eq (3) for some $\rho > 2L$.

(i) We now prove that $\mathbf{x}^*$ is also a global minimizer of Eq (3). For any feasible $\mathbf{x}$ and $\mathbf{v}$ that $\|\mathbf{x}\|_\infty \leq 1$, $\|\mathbf{v}\|_2^2 \leq n$, we derive the following inequalities:

$$
\begin{aligned}
& \mathcal{J}(\mathbf{x}, \mathbf{v}, \rho) \\
\geq\ & \min_{\|\mathbf{x}\|_\infty \leq 1,\ \|\mathbf{v}\|_2^2 \leq n,\ \mathbf{x} \in \Omega} f(\mathbf{x}) + \rho(n - \mathbf{x}^T\mathbf{v}) \\
=\ & \min_{\|\mathbf{x}\|_\infty \leq 1,\ \|\mathbf{v}\|_2^2 \leq n,\ \mathbf{x} \in \Omega} f(\mathbf{x}),\ s.t.\ \mathbf{x}^T\mathbf{v} = n \\
=\ & f(\mathbf{x}^*) + \rho(n - \mathbf{x}^{*T}\mathbf{v}^*) \\
=\ & \mathcal{J}(\mathbf{x}^*, \mathbf{v}^*, \rho)
\end{aligned}
$$

where the first equality holds due to the fact that the constraint $\mathbf{x}^T\mathbf{v} = n$ is satisfied at the local optimal solution when $\rho > 2L$ (see Lemma 7). Therefore, we conclude that

any optimal solution of Eq (2) is also an optimal solution of Eq (3).

(ii) We now prove that $\mathbf{x}_\rho^*$ is also a global minimizer of Eq (2). For any feasible $\mathbf{x}$ and $\mathbf{v}$ that $\|\mathbf{x}\|_\infty \leq 1$, $\|\mathbf{v}\|_2^2 \leq n$, $\mathbf{x}^T\mathbf{v} = n$, $\mathbf{x} \in \Omega$, we naturally have the following inequalities:

$$
\begin{aligned}
& f(\mathbf{x}_\rho^*) - f(\mathbf{x}) \\
=\ & f(\mathbf{x}_\rho^*) + \rho(n - \mathbf{x}_\rho^{*T}\mathbf{v}_\rho^*) \\
& -f(\mathbf{x}) - \rho(n - \mathbf{x}^T\mathbf{v}) \\
=\ & \mathcal{J}_\rho(\mathbf{x}_\rho^*, \mathbf{v}_\rho^*) - \mathcal{J}_\rho(\mathbf{x}, \mathbf{v}) \\
\leq\ & 0
\end{aligned}
$$

where the first equality uses Lemma 7. Therefore, we conclude that any optimal solution of Eq (3) is also an optimal solution of Eq (2).

Finally, we conclude that when $\rho > 2L$, the biconvex optimization in Eq (3) has the same local and global minima with the original problem in Eq (2). $\qquad \square$

**Theorem 2.** *Convergence Rate and Asymptotic Monotone Property of Algorithm 1. Assume that $f(\cdot)$ is a L-Lipschitz continuous convex function on $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. Algorithm 1 will converge to the first-order KKT point in at most $\lceil (\ln(L\sqrt{2n}) - \ln(\epsilon\rho^0))/\ln\sigma \rceil$ outer iterations with the accuracy at least $n - \mathbf{x}^T\mathbf{v} \leq \epsilon$. Moreover, after $\langle \mathbf{x}, \mathbf{v} \rangle = n$ is obtained, the sequence of $\{f(\mathbf{x}^t)\}$ generated by Algorithm 1 is monotonically non-increasing.*

*Proof.* We denote $s$ and $t$ as the outer iterations counter and inner iteration counter in Algorithm 1, respectively.

(i) we now prove the convergence rate of Algorithm 1. Assume that Algorithm 1 takes $s$ outer iterations to converge. We denote $f'(\mathbf{x})$ as the sub-gradient of $f(\cdot)$ in $\mathbf{x}$. According the the $\mathbf{x}$-subproblem in Eq (22), if $\mathbf{x}^*$ solves Eq (22), then we have the following variational inequality [26]:

$$
\begin{aligned}
& \forall \mathbf{x} \in [-1, +1]^n \cap \Omega,\ \langle \mathbf{x} - \mathbf{x}^*, f'(\mathbf{x}^*) \rangle + \\
& \rho(n - \sqrt{n}\|\mathbf{x}\|_2) - \rho(n - \sqrt{n}\|\mathbf{x}^*\|_2) \geq 0
\end{aligned}
$$

Letting $\mathbf{x}$ be any feasible solution that $\mathbf{x} \in \{-1, +1\}^n \cap \Omega$, we have the following inequality:

$$
\begin{aligned}
& (n - \sqrt{n}\|\mathbf{x}^*\|_2) \\
\leq\ & (n - \sqrt{n}\|\mathbf{x}\|_2) + \tfrac{1}{\rho}\langle \mathbf{x} - \mathbf{x}^*, f'(\mathbf{x}^*) \rangle \\
\leq\ & \tfrac{1}{\rho}\|\mathbf{x} - \mathbf{x}^*\|_2 \cdot \|f'(\mathbf{x}^*)\|_2 \\
\leq\ & L\sqrt{2n}/\rho \quad (24)
\end{aligned}
$$

where the second inequality is due to the Cauchy-Schwarz Inequality, the third inequality is due to the fact that $\|\mathbf{x} - \mathbf{y}\|_2 \leq \sqrt{2n}$, $\forall -\mathbf{1} \leq \mathbf{x}, \mathbf{y} \leq \mathbf{1}$ and the Lipschitz continuity of $f(\cdot)$ that $\|f'(\mathbf{x}^*)\|_2 \leq L$.

18

The inequality in Eq (24) implies that when $\rho^s \geq L\sqrt{2n}/\epsilon$, Algorithm 1 achieves accuracy at least $n - \sqrt{n}\|\mathbf{x}\|_2 \leq \epsilon$. Noticing that $\rho^s = \sigma^s \rho^0$, we have that $\epsilon$ accuracy will be achieved when

$$\sigma^s \rho^0 \geq \frac{L\sqrt{2n}}{\epsilon}$$
$$\Rightarrow \quad \sigma^s \geq \frac{L\sqrt{2n}}{\epsilon \rho^0}$$
$$\Rightarrow \quad s \geq (\ln(L\sqrt{2n}) - \ln(\epsilon \rho^0))/\ln \sigma$$

(ii) we now prove the asymptotic monotone property of Algorithm 1. We naturally derive the following inequalities:

$$
\begin{aligned}
&f(\mathbf{x}^{t+1}) - f(\mathbf{x}^t) \\
\leq\ & \rho(n - \langle \mathbf{x}^t, \mathbf{v}^t \rangle) - \rho(n - \langle \mathbf{x}^{t+1}, \mathbf{v}^t \rangle) \\
=\ & \rho\left(\langle \mathbf{x}^{t+1}, \mathbf{v}^t \rangle - \langle \mathbf{x}^t, \mathbf{v}^t \rangle\right) \\
\leq\ & \rho\left(\langle \mathbf{x}^{t+1}, \mathbf{v}^{t+1} \rangle - \langle \mathbf{x}^t, \mathbf{v}^t \rangle\right) \\
=\ & 0
\end{aligned}
$$

where the first step uses the fact that $f(\mathbf{x}^{t+1}) + \rho(n - \langle \mathbf{x}^{t+1}, \mathbf{v}^t \rangle) \leq f(\mathbf{x}^t) + \rho(n - \langle \mathbf{x}^t, \mathbf{v}^t \rangle)$ holds due to $\mathbf{x}^{t+1}$ is the optimal solution of Eq (6); the third step uses the fact $-\langle \mathbf{x}^{t+1}, \mathbf{v}^{t+1} \rangle \leq -\langle \mathbf{x}^{t+1}, \mathbf{v}^t \rangle$ holds due to $\mathbf{v}^{t+1}$ is the optimal solution of Eq (7); the last step uses $\langle \mathbf{x}, \mathbf{v} \rangle = n$. Note that the equality $\langle \mathbf{x}, \mathbf{v} \rangle = n$ together with the feasible set $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$, $\|\mathbf{v}\|_2^2 \leq n$ also implies that $\mathbf{x} \in \{-1, +1\}^n$. $\square$

The following lemma is useful in building the exactness property of the augmented Lagrangian function in Algorithm 2.

**Lemma 8.** *Consider the following optimization problem:*

$$(\mathbf{x}_\rho^*, \mathbf{v}_\rho^*) = \arg \min_{-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \|\mathbf{v}\|_2^2 \leq n, \mathbf{x} \in \Omega} \mathcal{L}(\mathbf{x}, \mathbf{v}, \rho). \quad (25)$$

*Assume that $f(\cdot)$ is a $L$-Lipschitz continuous convex function on $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. When $\rho > 2L$, $\langle \mathbf{x}_\rho^*, \mathbf{v}_\rho^* \rangle = n$ will be achieved for any local optimal solution of Eq (25).*

*Proof.* The proof of this theorem is based on Lemma 7. We observe that $n - \mathbf{x}^T \mathbf{v} = 0 \Leftrightarrow (n - \mathbf{x}^T \mathbf{v})^2 = 0$. We define $h(\mathbf{x}) \triangleq f(\mathbf{x}) + \frac{\alpha}{2}(n - \mathbf{x}^T \mathbf{v})^2$ and denote $L_h$ as the Lipschtz constant of $h(\cdot)$. We replace $f(\cdot)$ with $h(\cdot)$ in Lemma 7 and conclude that when $\rho > 2L_g$, we have $n - \mathbf{x}^T \mathbf{v} = 0$. Thus, the term $\frac{\alpha}{2}(n - \mathbf{x}^T \mathbf{v})^2$ in $h(\mathbf{x})$ reduces to zero and $L_h = L$. We conclude that when $\rho > 2L$, $\langle \mathbf{x}_\rho^*, \mathbf{v}_\rho^* \rangle = n$ will be achieved for any local optimal solution. $\square$

**Theorem 3.** *Exactness of the augmented Lagrangian Function. Assume that $f(\cdot)$ is a $L$-Lipschitz continuous convex function on $-\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}$. When $\rho > 2L$, the biconvex optimization problem $\min_{\mathbf{x}, \mathbf{v}} \mathcal{L}(\mathbf{x}, \mathbf{v}, \rho), s.t. -\mathbf{1} \leq \mathbf{x} \leq \mathbf{1}, \|\mathbf{v}\|_2^2 \leq n, \mathbf{x} \in \Omega$ in Eq (9) has the same local and global minima with the original problem in Eq (2).*

*Proof.* The proof of this theorem is similar to Theorem 1. We let $\mathbf{x}^*$ be any global minimizer of Eq (2) and $(\mathbf{x}_\rho^*, \mathbf{v}_\rho^*)$ be any global minimizer of Eq (3) for some $\rho > 2L$.

(i) We now prove that $\mathbf{x}^*$ is also a global minimizer of Eq (3). For any feasible $\mathbf{x}$ and $\mathbf{v}$ that $\|\mathbf{x}\|_\infty \leq 1$, $\|\mathbf{v}\|_2^2 \leq n$, $\mathbf{x} \in \Omega$, we derive the following inequalities:

$$
\begin{aligned}
&\mathcal{L}(\mathbf{x}, \mathbf{v}, \rho) \\
\geq\ & \min_{\|\mathbf{x}\|_\infty \leq 1, \|\mathbf{v}\|_2 \leq n} f(\mathbf{x}) + \rho(n - \mathbf{x}^T \mathbf{v}) + \frac{\alpha}{2}(n - \mathbf{x}^T \mathbf{v})^2 \\
=\ & \min_{\|\mathbf{x}\|_\infty \leq 1, \|\mathbf{v}\|_2 \leq n} f(\mathbf{x}), \ s.t. \ \mathbf{x}^T \mathbf{v} = n \\
=\ & f(\mathbf{x}^*) + \rho(n - \mathbf{x}^{*T} \mathbf{v}^*) + \frac{\alpha}{2}(n - \mathbf{x}^{*T} \mathbf{v}^*)^2 \\
=\ & \mathcal{L}(\mathbf{x}^*, \mathbf{v}^*, \rho)
\end{aligned}
$$

where the first equality holds due to the fact that the constraint $\mathbf{x}^T \mathbf{v} = n$ is satisfied at the local optimal solution when $\rho > 2L$ (see Lemma 7). Therefore, we conclude that any optimal solution of Eq (2) is also an optimal solution of Eq (3).

(ii) We now prove that $\mathbf{x}_\rho^*$ is also a global minimizer of Eq (2). For any feasible $\mathbf{x}$ and $\mathbf{v}$ that $\|\mathbf{x}\|_\infty \leq 1$, $\|\mathbf{v}\|_2^2 \leq n$, $\mathbf{x}^T \mathbf{v} = n$, $\mathbf{x} \in \Omega$, we naturally have the following inequalities:

$$
\begin{aligned}
&f(\mathbf{x}_\rho^*) - f(\mathbf{x}) \\
=\ & f(\mathbf{x}_\rho^*) + \rho(n - \mathbf{x}_\rho^{*T} \mathbf{v}_\rho^*) + \frac{\alpha}{2}(n - \mathbf{x}_\rho^{*T} \mathbf{v}_\rho^*)^2 \\
& -f(\mathbf{x}) - \rho(n - \mathbf{x}^T \mathbf{v}) - \frac{\alpha}{2}(n - \mathbf{x}^T \mathbf{v})^2 \\
=\ & \mathcal{L}(\mathbf{x}_\rho^*, \mathbf{v}_\rho^*, \rho) - \mathcal{L}(\mathbf{x}, \mathbf{v}, \rho) \\
\leq\ & 0
\end{aligned}
$$

Therefore, we conclude that any optimal solution of Eq (3) is also an optimal solution of Eq (2).

Finally, we conclude that when $\rho > 2L$, the biconvex optimization in Eq (3) has the same local and global minima with the original problem in Eq (2).

$\square$

## C. Solving the Rank-One Subproblem

This subsection describes a nearly closed-form solution for solving the rank-one subproblem which is involved in our alternating direction method. For general purpose, we consider the following optimization problem:

$$\min_{\|\mathbf{x}\|_2 \leq \beta} \frac{1}{2}\mathbf{x}^T(\gamma \mathbf{I} + \mathbf{b}\mathbf{b}^T)\mathbf{x} + \langle \mathbf{x}, \mathbf{c} \rangle \quad (26)$$

where $\gamma, \beta \in \mathbb{R}$, $\mathbf{b}, \mathbf{c} \in \mathbb{R}^n$ are given. We assume $\beta \neq 0$ or $\gamma \neq 0$. Clearly, Eq (26) is equivalent to the following minimax optimization problem:

$$\min_{\mathbf{x}} \max_{\theta \geq 0} \frac{1}{2}\mathbf{x}^T(\gamma \mathbf{I} + \mathbf{b}\mathbf{b}^T)\mathbf{x} + \langle \mathbf{x}, \mathbf{c} \rangle + \frac{\theta}{2}(\|\mathbf{x}\|_2^2 - \beta^2) \quad (27)$$

Setting the gradient respect of $\mathbf{x}$ to zero, we obtain:

$$\mathbf{x} = -(\gamma\mathbf{I} + \theta\mathbf{I} + \mathbf{b}\mathbf{b}^T)^{-1}\mathbf{c} \tag{28}$$

Putting this equality into Eq (27), we have the following minimization with respect to $\theta$:

$$\max_{\theta \geq 0} -\frac{1}{2}\mathbf{c}^T(\gamma\mathbf{I} + \theta\mathbf{I} + \mathbf{b}\mathbf{b}^T)^{-1}\mathbf{c} - \frac{1}{2}\theta\beta^2 \tag{29}$$

Using the well-known Sherman-Morrison inverse formula [22], Eq (29) reduces to the following optimization problem:

$$\max_{\theta \geq 0} \frac{1}{2}\left(\frac{t^2}{(\gamma + \theta)(\gamma + \theta + r)} - \frac{s}{\gamma + \theta} - s\theta\beta^2\right)$$

where $r = \mathbf{b}^T\mathbf{b}$, $s = \mathbf{c}^T\mathbf{c}$, $t = \mathbf{c}^T\mathbf{b}$. The optimal solution $\theta^*$ can be found using a simple one-dimensional bisection line search procedure. After that, the optimal solution $\mathbf{x}^*$ for the original optimization problem in Eq (26) can be recovered using Eq (28).

---

[22] $(\eta\mathbf{I} + \mathbf{b}\mathbf{b}^T)^{-1} = \frac{1}{\eta}\mathbf{I} - \frac{1}{\eta^2 + \eta\mathbf{b}^T\mathbf{b}}\mathbf{b}\mathbf{b}^T$, $\forall\eta > 0$