

# MAP Inference via $\ell_2$ -Sphere Linear Program Reformulation

Baoyuan Wu · Li Shen · Tong Zhang · Bernard Ghanem

the date of receipt and acceptance should be inserted later

**Abstract** Maximum a posteriori (MAP) inference is an important task for graphical models. Due to complex dependencies among variables in realistic models, finding an exact solution for MAP inference is often intractable. Thus, many approximation methods have been developed, among which the linear programming (LP) relaxation based methods show promising performance. However, one major drawback of LP relaxation is that it is possible to give fractional solutions. Instead of presenting a tighter relaxation, in this work we propose a continuous but equivalent reformulation of the original MAP inference problem, called LS-LP. We add the  $\ell_2$ -sphere constraint onto the original LP relaxation, leading to an intersected space with the local marginal polytope that is equivalent to the space of all valid integer label configurations. Thus, LS-LP is equivalent to the original MAP inference problem. We propose a perturbed alternating direction

method of multipliers (ADMM) algorithm to optimize the LS-LP problem, by adding a sufficiently small perturbation  $\epsilon$  onto the objective function and constraints. We prove that the perturbed ADMM algorithm globally converges to the  $\epsilon$ -Karush–Kuhn–Tucker ( $\epsilon$ -KKT) point of the LS-LP problem. The convergence rate will also be analyzed. Experiments on several benchmark datasets from Probabilistic Inference Challenge (PIC 2011) and OpenGM 2 show competitive performance of our proposed method against state-of-the-art MAP inference methods.

## 1 Introduction

Given the probability distribution of a graphical model, maximum a posteriori (MAP) inference aims to infer the most probable label configuration. MAP inference can be formulated as an integer linear program (ILP) [39]. However, due to the integer constraint, the exact optimization of ILP is intractable in many realistic problems. To tackle it, a popular approach is relaxing ILP to a continuous linear program over a local marginal polytope, *i.e.*,  $\mathcal{L}_G$  (defined in Section 3), called linear programming (LP) relaxation. The optimal solution to the LP relaxation will be obtained at the vertices of  $\mathcal{L}_G$ . It has been known [39] that all valid integer label configurations are at the vertices of  $\mathcal{L}_G$ , but not all vertices of  $\mathcal{L}_G$  are integer, while some are fractional. Since LP relaxation is likely to give fractional solutions, the rounding method must be adopted to generate integer solutions. To alleviate this issue, intense efforts have been made to design tighter relaxations (*e.g.*, high-order relaxation [36]) based on LP relaxation, such that the proportion of fractional vertices of  $\mathcal{L}_G$  can be reduced. However, the possibility of fractional solutions still exists. And, these tighter relaxations are often much more computationally expensive than the original LP relaxation. Moreover, there are also exact inference meth-

---

Baoyuan Wu was partially supported by Tencent AI Lab and King Abdullah University of Science and Technology (KAUST). Li Shen was supported by Tencent AI Lab. Bernard Ghanem was supported by the King Abdullah University of Science and Technology (KAUST) Office of Sponsored Research through the Visual Computing Center (VCC) funding. Tong Zhang was supported by the Hong Kong University of Science and Technology (HKUST). Li Shen is the corresponding author.

---

Baoyuan Wu  
Tencent AI Lab, Shenzhen 518000, China  
E-mail: wubaoyuan1987@gmail.com

Li Shen  
Tencent AI Lab, Shenzhen 518000, China  
E-mail: mathshenli@gmail.com

Tong Zhang  
Hong Kong University of Science and Technology, Hong Kong, China  
E-mail: tongzhang@tongzhang-ml.org

Bernard Ghanem  
King Abdullah University of Science and Technology, Thuwal 23955, Saudi Arabia  
E-mail: bernard.ghanem@kaust.edu.sa

ods, such as branch-and-bound [23] and cutting-plane [17], which utilize LP relaxation as sub-routines, leading to much higher computational cost than approximate methods.

Instead of proposing a new approximation with a tighter relaxation, we propose an exact reformulation of the original MAP inference problem. Specifically, we add a new constraint, called  $\ell_2$ -sphere [41], onto the original LP relaxation problem. It enforces that the solution  $\mathbf{x} \in \mathbb{R}^n$  should be on a  $\ell_2$ -sphere, *i.e.*,  $\|\mathbf{x} - \frac{1}{2}\|_2^2 = \frac{n}{4}$ . We can prove that the intersection between the  $\ell_2$ -sphere constraint and the local polytope  $\mathcal{L}_G$  is equivalent to the set of all possible label configurations of the original MAP inference problem, *i.e.*, the constraint space of the ILP problem. Thus, the proposed formulation, dubbed LS-LP, is an equivalent but continuous reformulation of the ILP formulation for MAP inference. Furthermore, inspired by [28] and [41], we adopt the ADMM algorithm [6], to not only separate the different constraints, but also decompose variables to allow parallel inference by exploiting the factor graph structure. Although the  $\ell_2$ -sphere constraint is non-convex, we prove that the ADMM algorithm for the LS-LP problem with a sufficiently small perturbation  $\epsilon$  will globally converges to the  $\epsilon$ -KKT [16,22] point of the original LS-LP problem. The obvious advantages of the proposed LS-LP formulation and the corresponding ADMM algorithm include: **1)** compared to other LP relaxation based methods, our method directly gives the valid integer label configuration, without any rounding techniques as post-processing; **2)** compared to the exact methods like branch-and-bound [23] and cutting-plane [17], our method optimizes one single continuous problem once, rather than multiple times. Experiments on benchmarks from Probabilistic Inference Challenge (PIC 2011) [7] and OpenGM 2 [14] verify the competitive performance of LS-LP against state-of-the-art MAP inference methods.

The main contributions of this work are three-fold. **1)** We propose a continuous but equivalent reformulation of the MAP inference problem. **2)** We present the ADMM algorithm for optimizing the perturbed LS-LP problem, which is proved to be globally convergent to the  $\epsilon$ -KKT point of the original LS-LP problem. The analysis of convergence rate is also presented. **3)** Experiments on benchmark datasets verify the competitive performance of our method compared to state-of-the-art MAP inference methods.

## 2 Related Work

As our method is closely related to LP relaxation based MAP inference methods, here we mainly review MAP inference methods of this category. For other categories of methods, such as message passing and move making, we refer the readers to [39] and [14] for more details. Although some off-the-shelf LP solvers can be used to optimize the LP relaxation problem, in many real-world applications the prob-

lem scale is too large to adopt these solvers. Hence, most methods focus on developing efficient algorithms to optimize the dual LP problem. Block coordinate descent methods [9,19] are fast, but they may converge to sub-optimal solutions. Sub-gradient based methods [20,15] can converge to global solutions, but their convergence is slow. Their common drawback is the non-smoothness of the dual objective function. To handle this difficulty, some smoothing methods have been developed. The Lagrangian relaxation [12] method uses the smooth log-sum-exp function to approximate the non-smooth max function in the dual objective. A proximal regularization [13] or an  $\ell_2$  regularization term [30] is added to the dual objective. Moreover, the steepest  $\epsilon$ -descent method proposed in [34] and [35] can accelerate the convergence of the standard sub-gradient based methods. Parallel MAP inference methods based on ADMM have also been developed to handle large-scale inference problems. For example, AD3 [28,27] and Bethe-ADMM [8] optimize the primal LP problem, while ADMM-dual [29] optimizes the dual LP problem. The common drawback of these methods is that they are likely to produce fractional solutions, since the underlying problem is merely a relaxation to the MAP inference problem.

Another direction is pursuing tighter relaxations, such as high-order consistency [36] and SDP relaxation [24]. But they are often more computationally expensive than LP relaxations. In contrast, the formulation of the proposed LS-LP is an exact reformulation of the original MAP inference problem, and the adopted ADMM algorithm can explicitly produce valid integer label configurations, without any rounding operation. In comparison with other expensive exact MAP inference methods (*e.g.*, Branch-and-Bound [23] and cutting plane [17]), LS-LP is very efficient owing to the resulting parallel inference, similar to other ADMM based methods.

Another related work is  $\ell_p$ -Box ADMM [41], which is a framework to optimize the general integer program. The proposed LS-LP is inspired by this framework, where the integer constraints are replaced by the intersection of two continuous constraints. However, **1)** LS-LP is specifically designed for MAP inference, as it replaces the valid integer configuration space (*e.g.*,  $\{(0,1), (1,0)\}$  for the variable with binary states), rather than the whole binary space (*e.g.*,  $\{(0,0), (0,1), (1,0), (1,1)\}$ ) as did in  $\ell_p$ -Box ADMM. **2)** LS-LP is tightly combined with LP relaxation, and the ADMM algorithm decomposes the problem into multiple simple sub-problems by utilizing the structure of the factor graph, which allows parallel inference for any type of inference problems (*e.g.*, multiple variable states and high-order factors). In contrast,  $\ell_p$ -Box ADMM does not assume any special property for the objective function, and it optimizes all variable nodes in one sub-problem. Especially for large-scale models, the sub-problem involved in  $\ell_p$ -Box ADMM will be very cost. **3)** As LP relaxation is parameterized according to the factor

graph, any type of graphical models (*e.g.*, directed models, high-order potentials, asymmetric potentials) can be naturally handled by LS-LP. In contrast,  $\ell_p$ -Box ADMM needs to transform the inference objective based on MRF models to some simple forms (*e.g.*, binary quadratic program (BQP)). However, the transformation is non-trivial in some cases. For example, if there are high-order potentials, the graphical model is difficult to input into a BQP problem.

### 3 Background

#### 3.1 Factor Graph

Denote  $\mathbf{G} = \{\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_N\}$  as a set of  $N$  random variables in a discrete space  $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_N$ , where  $\mathcal{X}_i = \{0, \dots, r_i - 1\}$  with  $r_i = |\mathcal{X}_i|$  being the number of possible states of  $\mathbf{g}_i$ . The joint probability of  $\mathbf{G}$  is formulated based on a factor graph  $G$  [18],

$$P(\mathbf{G}) \propto \exp\left(\sum_{i \in V} \theta_i(\mathbf{g}_i) + \sum_{\alpha \in F} \theta_\alpha(\mathbf{g}_\alpha)\right), \quad (1)$$

where  $G = (V, F, E)$  with  $V = \{1, \dots, N\}$  being the node set of variables,  $F$  being the node set of factors, as well as the edge set  $E \subseteq V \times F$  linking the variable and factor nodes. A simple MRF model and its factor graph are shown in Fig. 1(a,b). We refer the readers to [18] for the detailed definition of the factor graph.  $\mathbf{g}_\alpha$  indicates the label configuration of the factor  $\alpha$ , and its state will be determined according to the states of connected variable nodes, *i.e.*,  $\{\mathbf{g}_i | i \in \mathcal{N}_\alpha\}$ , with  $\mathcal{N}_\alpha$  being the set of neighborhood variable nodes of the factor  $\alpha$ .  $\theta_i(\cdot)$  denotes the unary log potential (logPot) function, while  $\theta_\alpha(\cdot)$  indicates the factor logPot function.

#### 3.2 MAP Inference as Linear Program

Given  $P(\mathbf{G})$ , an important task is to find the most probable label configuration of  $\mathbf{G}$ , referred to as MAP inference,

$$\text{MAP}(\theta) = \max_{\mathbf{G} \in \mathcal{X}} \sum_{i \in V} \theta_i(\mathbf{g}_i) + \sum_{\alpha \in F} \theta_\alpha(\mathbf{g}_\alpha). \quad (2)$$

Eq. (2) can be reformulated as the integer linear program (ILP) [39],

$$\text{ILP}(\theta) = \max_{\boldsymbol{\mu}} \sum_{i \in V} \theta_i^\top \boldsymbol{\mu}_i + \sum_{\alpha \in F} \theta_\alpha^\top \boldsymbol{\mu}_\alpha = \max_{\boldsymbol{\mu}} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle, \quad (3)$$

s.t.  $\boldsymbol{\mu} \in \mathcal{L}_G \cap \{0, 1\}^{|\boldsymbol{\mu}|}$ ,

where  $\boldsymbol{\theta} = (\dots; \theta_i; \dots; \theta_\alpha; \dots)$ ,  $i \in V, \alpha \in F$  denotes the log potential (logPot) vector, derived from  $\theta_i(\mathbf{g}_i)$  and  $\theta_\alpha(\mathbf{g}_\alpha)$ .  $\boldsymbol{\mu} = [\boldsymbol{\mu}_V; \boldsymbol{\mu}_F]$ , where  $\boldsymbol{\mu}_V = [\boldsymbol{\mu}_1; \dots; \boldsymbol{\mu}_{|V|}]$  and  $\boldsymbol{\mu}_F = [\boldsymbol{\mu}_1; \dots; \boldsymbol{\mu}_{|F|}]$ .  $\boldsymbol{\mu}_i \in \{0, 1\}^{|\mathcal{X}_i|}$  indicates the label

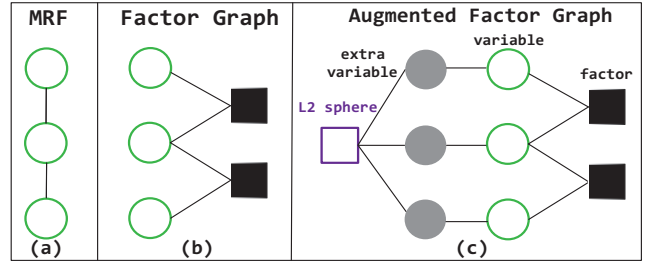


Fig. 1: An example of (a) MRF, (b) factor graph corresponding to LP and (c) augmented factor graph corresponding to LS-LP.

vector corresponding to  $\mathbf{g}_i$ : if the state of  $\mathbf{g}_i$  is  $t$ , then  $\boldsymbol{\mu}_i(t) = 1$ , while all other entries are 0. Similarly,  $\boldsymbol{\mu}_\alpha \in \{0, 1\}^{|\mathcal{X}_\alpha|}$  indicates the label vector corresponding to  $\mathbf{g}_\alpha$ . The local marginal polytope is defined as follows,

$$\mathcal{L}_G = \left\{ \boldsymbol{\mu} \mid \boldsymbol{\mu}_\alpha \in \Delta^{|\mathcal{X}_\alpha|}, \forall \alpha \in F; \right. \\ \left. \boldsymbol{\mu}_i = \mathbf{M}_{i\alpha} \boldsymbol{\mu}_\alpha, \forall (i, \alpha) \in E \right\}, \quad (4)$$

with  $\Delta^{|\mathbf{a}|} = \{\mathbf{a} \mid \mathbf{1}^\top \mathbf{a} = 1, \mathbf{a} \geq \mathbf{0}\}$  being the probability simplex, and the second constraint ensures the local consistency between  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\mu}_\alpha$ .  $\mathbf{M}_{i\alpha} \in \{0, 1\}^{|\mathcal{X}_i| \times |\mathcal{X}_\alpha|}$  of the local consistency constraint included in  $\mathcal{L}_G$  is defined as: the entry of  $\mathbf{M}_{i\alpha}$  is 1 if  $\mathbf{g}_\alpha \sim \mathbf{g}_i$ , where  $\mathbf{g}_\alpha \sim \mathbf{g}_i$  indicates the state of  $\mathbf{g}_i$  and the state of the corresponding element in  $\mathbf{g}_\alpha$  are the same; otherwise, the entry is 0. For example, we consider a binary-state variable node  $\boldsymbol{\mu}_i \in \{0, 1\}^2$  and a pairwise factor node  $\boldsymbol{\mu}_\alpha \in \{0, 1\}^4$  connected to two variable nodes (the variable node  $i$  is the first). The first entry of  $\boldsymbol{\mu}_i$  indicates the score of choosing state 0, while the second entry corresponds to that of choosing state 1. The four entries of  $\boldsymbol{\mu}_\alpha$  indicate the scores of four label configurations of two connected variables, *i.e.*,  $(0, 0), (0, 1), (1, 0), (1, 1)$ . In this case,  $\mathbf{M}_{i\alpha} = [1, 1, 0, 0; 0, 0, 1, 1]$ .

Moreover, Eq. (2) can also be rewritten as

$$\text{MAP}(\theta) = \max_{\boldsymbol{\mu} \in \mathcal{M}_G} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle, \quad (5)$$

where the marginal polytope is defined as follows,

$$\mathcal{M}_G = \{\boldsymbol{\mu} \mid \exists P(X), \text{ such that } \boldsymbol{\mu}_i, \boldsymbol{\mu}_\alpha \in \mathcal{L}_G\}. \quad (6)$$

Solving  $\text{MAP}(\theta)$  is difficult (NP-hard in general), especially for large scale problems. Instead, the approximation over  $\mathcal{L}_G$  is widely adopted, as follows:

$$\text{LP}(\theta) = \max_{\boldsymbol{\mu} \in \mathcal{L}_G} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle \geq \text{ILP}(\theta) = \text{MAP}(\theta), \quad (7)$$

which is called LP relaxation. Note that here  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\mu}_\alpha$  are continuous variables, and they are considered as local marginals of  $\mathbf{g}_i$  and  $\mathbf{g}_\alpha$ , respectively.

According to [39], the characteristics of  $\text{LP}(\theta)$ ,  $\text{MAP}(\theta)$  and their relationships are briefly summarized in Lemma 1.

**Lemma 1** [39] *The relationship between  $\mathcal{M}_G$  and  $\mathcal{L}_G$ , and that between  $\text{MAP}(\theta)$  and  $\text{LP}(\theta)$  are as follows.*

- $\mathcal{M}_G \subseteq \mathcal{L}_G$ ;
- $\text{MAP}(\theta) \leq \text{LP}(\theta)$ ;
- All vertices of  $\mathcal{M}_G$  are integer, while  $\mathcal{L}_G$  includes both integer and fractional vertices. And the set of integer vertices of  $\mathcal{L}_G$  is same with the set of the vertices of  $\mathcal{M}_G$ . All non-vertices in  $\mathcal{M}_G$  and  $\mathcal{L}_G$  are fractional points.
- Since both  $\mathcal{M}_G$  and  $\mathcal{L}_G$  are convex polytopes, the global solutions of  $\text{MAP}(\theta)$  and  $\text{LP}(\theta)$  will be on the vertices of  $\mathcal{M}_G$  and  $\mathcal{L}_G$ , respectively.
- The global solution  $\mu^*$  of  $\text{LP}(\theta)$  can be fractional or integer. If it is integer, then it is also the global solution of  $\text{MAP}(\theta)$ .

### 3.3 Kurdyka-Lojasiewicz Inequality

The Kurdyka-Lojasiewicz inequality was firstly proposed in [26], and it has been widely used in many recent works [2, 40, 25] for the convergence analysis of non-convex problems. Since it will also be used in the later convergence analysis of our algorithm, it is firstly produced here, as shown in Definition 1.

**Definition 1** [2] A function  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup +\infty$  is said to have the Kurdyka-Lojasiewicz (KL) property at  $x^* \in \text{dom}(\partial f)$  ( $\text{dom}(\cdot)$  denotes the domain of function,  $\partial$  indicates the sub-gradient operator), if the following two conditions hold

- there exist a constant  $\eta \in (0, +\infty]$ , a neighborhood  $\mathcal{V}_{x^*}$  of  $x^*$ , as well as a continuous concave function  $\varphi : [0, \eta) \rightarrow \mathbb{R}_+$ , with  $\varphi(0) = 0$  and  $\varphi$  is differentiable on  $(0, \eta)$  with positive derivatives.
- $\forall x \in \mathcal{V}_{x^*}$  satisfying  $f(x^*) < f(x) < f(x^*) + \eta$ , the Kurdyka-Lojasiewicz inequality holds

$$\varphi'(f(x) - f(x^*)) \text{dist}(0, \partial f(x)) \geq 1. \quad (8)$$

**Remark.** According to [2, 4, 5], if  $f$  is semi-algebraic, then it satisfies the KL property with  $\varphi(s) = cs^{1-p}$ , where  $p \in [0, 1)$  and  $c > 0$  are constants. This point will be used in later analysis of convergence.

## 4 MAP Inference via $\ell_2$ -sphere Linear Program Reformulation

### 4.1 Equivalent Reformulation

Firstly, we introduce the  $\ell_2$ -sphere constraint [41],

$$\mathcal{S} = \left\{ \mathbf{x} \in \mathbb{R}^n \mid \left\| \mathbf{x} - \frac{1}{2} \mathbf{1} \right\|_2^2 = \frac{n}{4} \right\}. \quad (9)$$

Note that  $\mathcal{S}$  is defined with respect to the vector  $\mathbf{x}$ , rather than individual scalars  $x_i, i = 1, \dots, n$ . We propose to add the  $\ell_2$ -sphere constraint onto the variable nodes  $\mu_V$ . Combining this with LP relaxation (see Eq. (7)), we propose a new formulation for MAP inference,

$$\text{LS-LP}(\theta) = \max_{\mu} \langle \theta, \mu \rangle, \quad \text{s.t. } \mu \in \mathcal{L}_G, \mu_V \in \mathcal{S}. \quad (10)$$

Due to the non-convex constraint  $\mathcal{S}$ , it is no longer a linear program. However, to emphasize its relationship to LP relaxation, we still denote it as a  $\ell_2$ -sphere constrained linear program (LS-LP) reformulation. More importantly, as shown in Proposition 1, LS-LP is equivalent to the original MAP inference problem, rather than a relaxation as in LP. Inspired by the constraint separation in  $\ell_p$ -Box ADMM [41], we introduce the extra variable  $\mathbf{v}$  to reformulate (10) as

$$\text{LS-LP}(\theta) = \max_{\mu, \mathbf{v}} \langle \theta, \mu \rangle = \min_{\mu, \mathbf{v}} \langle -\theta, \mu \rangle, \quad (11)$$

$$\text{s.t. } \mu \in \mathcal{L}_G, \mathbf{v} \in \mathcal{S}, \mu_i = v_i, \forall i \in V,$$

where  $\mathbf{v} = [v_1; \dots; v_i; \dots; v_{|V|}]$ ,  $i \in V$  is the concatenated vector of all extra variable nodes. The combination of the original factor graph and these extra variable nodes is referred to as *augmented factor graph* (AFG). An example of AFG corresponding to Problem (11) is shown in Figure 1(c). The gray circles correspond to extra variables  $\mathbf{v}$ , and connections to the purple box indicate that  $\mathbf{v} \in \mathcal{S}$ . Note that AFG does not satisfy the definition of the standard factor graph, where connections only exist between variables nodes and factor nodes. However, AFG provides a clear picture of the structure of LS-LP and the node relationships. The proposed LS-LP problem is equivalent to the original MAP inference problem, as shown in Proposition 1. It means that the global solutions of this two problems are equivalent.

**Lemma 2** *The following constraint spaces are equivalent,*

$$\begin{aligned} \mathcal{C}_1 &= \{ \mu \mid \mu \in \mathcal{L}_G \cap \{0, 1\}^{|\mu|} \} \\ &\equiv \mathcal{C}_2 = \{ \mu \mid \mu \in \mathcal{L}_G \text{ and } \mu_V \in \mathcal{S} \} \\ &\equiv \mathcal{C}_3 = \{ \mu \mid \mu \in \mathcal{M}_G \cap \{0, 1\}^{|\mu|} \}. \end{aligned} \quad (12)$$

*Proof* We start from  $\mathcal{C}_2$ , where we have

$$\mu_V \in \mathcal{S} \iff \sum_{i \in V} \left\| \mu_i - \frac{1}{2} \right\|_2^2 = \frac{\sum_{i \in V} |\mathcal{X}_i|}{4}. \quad (13)$$

Besides, the following relations hold

$$\begin{aligned} \mu \in \mathcal{L}_G &\iff \mu_\alpha \in \Delta^{|\mathcal{X}_\alpha|} \text{ and } \mu_i = \mathbf{M}_{i\alpha} \mu_\alpha \\ \Rightarrow \mu_i \in [0, 1]^{|\mathcal{X}_i|} &\Rightarrow \left\| \mu_i - \frac{1}{2} \right\|_2^2 \leq \frac{|\mathcal{X}_i|}{4}, \end{aligned} \quad (14)$$

$\forall i \in V, \forall (i, \alpha) \in E$ . The equation in the last relation holds if and only if  $\mu_i \in \{0, 1\}^{|\mathcal{X}_i|}$ . Combining with (13), we conclude that  $\mu_i \in \{0, 1\}$  holds  $\forall i \in V$ . Consequently, utilizing the local consistency constraint  $\mu_i = \mathbf{M}_{i\alpha}\mu_\alpha$ , we obtain that  $\mu_\alpha \in \{0, 1\}$  also holds  $\forall \alpha \in F$ . Thus, we have  $\mu \in \{0, 1\}^{|\mu|}$ . Then, the relation  $\mathcal{C}_1 \equiv \mathcal{C}_2$  is proved.

Besides, as shown in Lemma 1, the set of integer vertices of  $\mathcal{L}_G$  is same with the one of  $\mathcal{M}_G$ , and all non-vertices in  $\mathcal{M}_G$  and  $\mathcal{L}_G$  are fractional points. Thus, it is easy to know  $\mathcal{C}_1 \equiv \mathcal{C}_3$ . Hence the proof is finished.

**Theorem 1** Utilizing Lemma 2, the aforementioned MAP inference problems have the following relationships,

$$LS-LP(\theta) = ILP(\theta) = MAP(\theta) \leq LP(\theta). \quad (15)$$

*Proof* According to Lemma 1.3 and 1.4, as well as  $\mathcal{C}_2 \equiv \mathcal{C}_3$  in Lemma 2 (see Eq. (12)), we have

$$\max_{\mu \in \mathcal{M}_G} \langle \theta, \mu \rangle = \max_{\mu \in \mathcal{C}_3} \langle \theta, \mu \rangle = \max_{\mu \in \mathcal{C}_2} \langle \theta, \mu \rangle \quad (16)$$

$$\iff MAP(\theta) = LS-LP(\theta). \quad (17)$$

Combining with  $MAP(\theta) = ILP(\theta) \leq LP(\theta)$  (see Eq. (7)), the proof is finished.

## 4.2 A General Form and KKT Conditions

For clarity, we firstly simplify the notations and formulations in Eq. (11) to the general shape,

$$LS-LP(\theta) = \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + h(\mathbf{y}), \text{ s.t. } \mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{y}. \quad (18)$$

Our illustration for (18) is separated into three parts, as follows:

- Variables.**  $\mathbf{x} = [\mu_1; \dots; \mu_{|V|}] \in \mathbb{R}^{\sum_i |\mathcal{X}_i|}$ , and it concatenates all variable nodes  $\mu_V$ .  $\mathbf{y} = [\mathbf{y}_1; \dots; \mathbf{y}_{|V|}]$  with  $\mathbf{y}_i = [\mathbf{v}_i; \mu_{\alpha_{i,1}}; \dots; \mu_{\alpha_{i,|\mathcal{N}_i|}}] \in \mathbb{R}^{|\mathcal{X}_i| + \sum_{\alpha} |\mathcal{X}_\alpha|}$ .  $\mathbf{y}$  concatenates all factor nodes  $\mu_V$  and the extra variable nodes  $\mathbf{v}$ ;  $\mathbf{y}_i$  concatenates the factor nodes and the extra variable node connected to the  $i$ -th variable node  $\mu_i$ .  $\mathcal{N}_i$  indicates the set of neighborhood factor nodes connected to the  $i$ -th variable node; the subscript  $\alpha_{i,j}$  indicates the  $j$ -th factor connected to the  $i$ -th variable, with  $i \in V$  and  $j \in \mathcal{N}_i$ .
- Objective functions.**  $f(\mathbf{x}) = \mathbf{w}_x^\top \mathbf{x}$  with  $\mathbf{w}_x = -[\theta_1; \dots; \theta_{|V|}]$ .  $h(\mathbf{y}) = g(\mathbf{y}) + \mathbf{w}_y^\top \mathbf{y}$ , with  $\mathbf{w}_y = [\mathbf{w}_1; \dots; \mathbf{w}_{|V|}]$  with  $\mathbf{w}_i = -[\mathbf{0}; \frac{1}{|\mathcal{N}_{\alpha_{i,1}}|} \theta_{\alpha_{i,1}}; \dots; \frac{1}{|\mathcal{N}_{\alpha_{i,|\mathcal{N}_i|}}|} \theta_{\alpha_{i,|\mathcal{N}_i|}}]$ , and  $\mathcal{N}_\alpha = \{i \mid (i, \alpha) \in E\}$  being the set of neighborhood variable nodes connected to the  $\alpha$ -th factor.  $g(\mathbf{y}) = \mathbb{I}(\mathbf{v} \in \mathcal{S}) + \sum_{\alpha \in F} \mathbb{I}(\mu_\alpha \in \Delta^{|\mathcal{X}_\alpha|})$ , with  $\mathbb{I}(a)$  being the indicator function:  $\mathbb{I}(a) = 0$  if  $a$  is true, otherwise  $\mathbb{I}(a) = \infty$ .

## Algorithm 1 The perturbed ADMM algorithm

**Input:** The initialization  $\mathbf{y}^0, \hat{\mathbf{x}}^0, \lambda^0$ , the perturbation  $\epsilon$ , the hyper-parameter  $\rho$

1: **for**  $k = 0$  to  $K$  **do**:

2:     Update  $\mathbf{y}^{k+1}$  as follows (see Section 5.1 for details)

$$\mathbf{y}^{k+1} = \underset{\mathbf{y}}{\operatorname{argmin}} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}^k, \lambda^k) \quad (21)$$

3:     Update  $\hat{\mathbf{x}}^{k+1}$  as follows (Section 5.2 for details)

$$\hat{\mathbf{x}}^{k+1} = \underset{\hat{\mathbf{x}}}{\operatorname{argmin}} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}, \lambda^k) \quad (22)$$

4:     Update  $\lambda^{k+1}$  (see Section 5.3 for details)

$$\lambda^{k+1} = \lambda^k + \rho(\hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}) \quad (23)$$

5:     Check stopping criterion, as shown in Section 5.4

6: **end for**

7: **return**  $\mathbf{y}^*, \mathbf{x}^*$

3. **Constraint matrices.** The constraint matrix  $\mathbf{A} = \operatorname{diag}(\mathbf{A}_1, \dots, \mathbf{A}_i, \dots, \mathbf{A}_{|V|})$  with  $\mathbf{A}_i = [\mathbf{I}_{|\mathcal{X}_i|}; \dots; \mathbf{I}_{|\mathcal{X}_i|}] \in \{0, 1\}^{(|\mathcal{N}_i|+1)|\mathcal{X}_i| \times |\mathcal{X}_i|}$ .  $\mathbf{B} = \operatorname{diag}(\mathbf{B}_1, \dots, \mathbf{B}_i, \dots, \mathbf{B}_{|V|})$ , with  $\mathbf{B}_i = \operatorname{diag}(\mathbf{I}_{|\mathcal{X}_i|}, \mathbf{M}_{i, \alpha_{i,1}}, \dots, \mathbf{M}_{i, \alpha_{i,|\mathcal{N}_i|}})$ .  $\mathbf{A}$  summarizes all constraints on  $\mu_V$ , while  $\mathbf{B}$  collects all constraints on  $\mu_F$  and  $\mathbf{v}$ .

Note that Problem (18) has a clear structure with two groups of variables, corresponding the augmented factor graph (see Fig. 1(c)).

**Definition 2** The solution  $(\mathbf{x}^*, \mathbf{y}^*)$  of the LS-LP problem (18) is said to be the KKT point if the following conditions are satisfied:

$$\mathbf{B}^\top \lambda^* \in \partial h(\mathbf{y}^*), \quad \nabla f(\mathbf{x}^*) = -\mathbf{A}^\top \lambda^*, \quad \mathbf{A}\mathbf{x}^* = \mathbf{B}\mathbf{y}^*, \quad (19)$$

where  $\lambda^*$  denotes the Lagrangian multiplier;  $\partial h$  indicates the sub-gradient of  $h$ , while  $\nabla f$  represents the gradient of  $f$ . Moreover,  $(\mathbf{x}^*, \mathbf{y}^*)$  is considered as the  $\epsilon$ -KKT point if the following conditions hold:

$$\begin{aligned} \operatorname{dist}(\mathbf{B}^\top \lambda^*, \partial h(\mathbf{y}^*)) &\leq O(\epsilon), \quad \|\nabla f(\mathbf{x}^*) + \mathbf{A}^\top \lambda^*\| \leq O(\epsilon), \\ \|\mathbf{A}\mathbf{x}^* - \mathbf{B}\mathbf{y}^*\| &\leq O(\epsilon). \end{aligned} \quad (20)$$

## 5 Perturbed ADMM Algorithm for LS-LP

We propose a perturbed ADMM algorithm to optimize the following perturbed augmented Lagrangian function,

$$\begin{aligned} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}, \lambda) &= \hat{f}(\hat{\mathbf{x}}) + h(\mathbf{y}) + \lambda^\top (\hat{\mathbf{A}}\hat{\mathbf{x}} - \mathbf{B}\mathbf{y}) \\ &\quad + \frac{\rho}{2} \|\hat{\mathbf{A}}\hat{\mathbf{x}} - \mathbf{B}\mathbf{y}\|_2^2, \end{aligned} \quad (24)$$

where  $\hat{\mathbf{A}} = [\mathbf{A}, \epsilon \mathbf{I}]$  with a sufficiently small constant  $\epsilon > 0$ , then  $\hat{\mathbf{A}}$  is full row rank.  $\hat{\mathbf{x}} = [\mathbf{x}; \bar{\mathbf{x}}]$ , with  $\bar{\mathbf{x}} = [\bar{\mathbf{x}}_1; \dots; \bar{\mathbf{x}}_{|V|}] \in$

$\mathbb{R}^{\sum_i^V (|\mathcal{N}_i|+1)|\mathcal{X}_i|}$  and  $\bar{\mathbf{x}}_i = [\boldsymbol{\mu}_i; \dots; \boldsymbol{\mu}_i] \in \mathbb{R}^{(|\mathcal{N}_i|+1)|\mathcal{X}_i|}$ .  $\hat{f}(\hat{\mathbf{x}}) = f(\mathbf{x}) + \frac{1}{2}\epsilon\hat{\mathbf{x}}^\top\hat{\mathbf{x}}$ . Note that both  $\hat{\mathbf{A}}$  and  $\mathbf{B}$  are full row rank, and the second-order gradient  $\nabla^2\hat{f}(\hat{\mathbf{x}}) = \epsilon\mathbf{I}$  is bounded. These properties will play key roles in our later analysis of convergence.

Following the conventional ADMM algorithm, the solution to the LS-LP problem (18) can be obtained through optimizing the following sub-problems based on (24) iteratively. The general structure of the algorithm is summarized in Algorithm 1.

### 5.1 Sub-Problem w.r.t. $\mathbf{y}$ in LS-LP Problem

Given  $\hat{\mathbf{x}}^k$  and  $\boldsymbol{\lambda}^k$ ,  $\mathbf{y}^{k+1}$  can be updated by solving the sub-problem (21) (see Algorithm 1). According to the definitions of  $\hat{\mathbf{A}}$ ,  $\hat{\mathbf{x}}$ ,  $\mathbf{B}$ ,  $\mathbf{y}$ , this problem can be further separated to the following two independent sub-problems, which can be solved in parallel.

**Update  $\mathbf{v}^{k+1}$ :**

$$\min_{\mathbf{v} \in \mathcal{S}} \sum_{i \in V} \left[ -(\boldsymbol{\lambda}_i^k)^\top \mathbf{v}_i + \frac{\rho_i}{2} \|(1+\epsilon)\boldsymbol{\mu}_i^k - \mathbf{v}_i\|_2^2 \right]. \quad (25)$$

It has a closed form solution as follows

$$\mathbf{v}^{k+1} = \mathcal{P}_{\mathcal{S}}(\bar{\mathbf{v}}^{k+1}), \quad (26)$$

where  $\bar{\mathbf{v}}^{k+1} = [\mathbf{v}_1^{k+1}; \dots; \mathbf{v}_{|V|}^{k+1}]$  with  $\mathbf{v}_i^{k+1} = (1+\epsilon)\boldsymbol{\mu}_i^k + \frac{1}{\rho_i}\boldsymbol{\lambda}_i^k$ .  $\mathcal{P}_{\mathcal{S}}(\cdot)$  is the projection onto  $\mathcal{S}$ :  $\mathcal{P}_{\mathcal{S}}(\mathbf{a}) = \frac{n^{1/2}}{\|\bar{\mathbf{a}}\|_2} \times \frac{\bar{\mathbf{a}}}{\|\bar{\mathbf{a}}\|_2} + \frac{1}{2}\mathbf{1}_n$ , with  $\bar{\mathbf{a}} = \mathbf{a} - \frac{1}{2}\mathbf{1}_n$  and  $n$  being the dimension of  $\mathbf{a}$ . As demonstrated in [41], this projected solution is the optimal solution to (25).

**Update  $\boldsymbol{\mu}_\alpha^{k+1}$ :** The sub-problems w.r.t.  $\{\boldsymbol{\mu}_\alpha\}_{\alpha \in F}$  can be run in parallel  $\forall \alpha \in F$ ,

$$\min_{\boldsymbol{\mu}_\alpha \in \Delta^{|\mathcal{X}_\alpha|}} -\boldsymbol{\theta}_\alpha^\top \boldsymbol{\mu}_\alpha + \sum_{i \in \mathcal{N}_\alpha} \left[ \frac{\rho_{i\alpha}}{2} \|(1+\epsilon)\boldsymbol{\mu}_i^k - \mathbf{M}_{i\alpha}\boldsymbol{\mu}_\alpha\|_2^2 - (\boldsymbol{\lambda}_{i\alpha}^k)^\top \mathbf{M}_{i\alpha}\boldsymbol{\mu}_\alpha \right]. \quad (27)$$

It is easy to know that Problem (27) is convex, as  $\mathbf{M}_{i\alpha}^\top \mathbf{M}_{i\alpha}$  is positive semi-definite and  $\Delta^{|\mathcal{X}_\alpha|}$  is a convex set. Any off-the-shelf QP solver can be adopted to solve (27). In experiments, we adopt the active-set algorithm implemented by a publicly-available toolbox called Quadratic Programming in C (QPC)<sup>1</sup>, which is written in C and can be called from MATLAB.

### 5.2 Sub-Problem w.r.t. $\hat{\mathbf{x}}$ in LS-LP Problem

Given  $\mathbf{y}^{k+1}$  and  $\boldsymbol{\lambda}^k$ ,  $\hat{\mathbf{x}}^{k+1}$  can be updated by solving the sub-problem (22) (see Algorithm 1). According to the definition of  $\hat{\mathbf{x}}$ , this problem can be separated to  $|V|$  independent sub-problems w.r.t.  $\{\boldsymbol{\mu}_i\}_{i \in V}$ , as follows:

$$\begin{aligned} \min_{\boldsymbol{\mu}_i} & (\boldsymbol{\lambda}_i^k - \boldsymbol{\theta}_i)^\top \boldsymbol{\mu}_i + \frac{\epsilon(|\mathcal{N}_i|+2)}{2} \boldsymbol{\mu}_i^\top \boldsymbol{\mu}_i + \sum_{\alpha \in \mathcal{N}_i} \left[ \right. \\ & \left. (1+\epsilon)(\boldsymbol{\lambda}_{i\alpha}^k)^\top \boldsymbol{\mu}_i + \frac{\rho_{i\alpha}}{2} \|(1+\epsilon)\boldsymbol{\mu}_i - \mathbf{M}_{i\alpha}\boldsymbol{\mu}_\alpha^{k+1}\|_2^2 \right] \\ & + \frac{\rho_i}{2} \|(1+\epsilon)\boldsymbol{\mu}_i - \mathbf{v}_i\|_2^2 \\ & = (1+\epsilon) \left[ \sum_{\alpha \in \mathcal{N}_i} (\boldsymbol{\lambda}_{i\alpha}^k - \rho_{i\alpha}\mathbf{M}_{i\alpha}\boldsymbol{\mu}_\alpha^{k+1}) - \rho_i\mathbf{v}_i^{k+1} - \boldsymbol{\theta}_i \right. \\ & \left. + \boldsymbol{\lambda}_i^k \right]^\top \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^\top \mathbf{Q}\boldsymbol{\mu}_i + \text{const}, \end{aligned} \quad (29)$$

where  $\mathbf{Q} = \frac{1}{2}[\epsilon(|\mathcal{N}_i|+2) + \rho_i(1+\epsilon)^2 + \sum_{\alpha \in \mathcal{N}_i} \rho_{i\alpha}(1+\epsilon)^2] \cdot \mathbf{I}$ . The above sub-problem can be further simplified to  $\min_{\boldsymbol{\mu}_i} \mathbf{a} \cdot \boldsymbol{\mu}_i^\top \boldsymbol{\mu}_i + \mathbf{b}^\top \boldsymbol{\mu}_i$ , where  $\mathbf{a}$  and  $\mathbf{b}$  can be easily derived from the above equation. Its close-form solution is obtained by setting its gradient to  $\mathbf{0}$ , i.e.,  $\boldsymbol{\mu}_i^{k+1} = \frac{\mathbf{b}}{2\mathbf{a}}$ .

### 5.3 Update $\boldsymbol{\lambda}$ in LS-LP Problem

Given  $\mathbf{y}^{k+1}$  and  $\hat{\mathbf{x}}^{k+1}$ ,  $\boldsymbol{\lambda}^{k+1}$  is updated using (23) (see Algorithm 1). Similarly, it can be separately to  $|V| + |E|$  independent sub-problems, as follows

$$\boldsymbol{\lambda}_i^{k+1} = \boldsymbol{\lambda}_i^k + \rho_i[(1+\epsilon)\boldsymbol{\mu}_i^{k+1} - \mathbf{v}_i^{k+1}], \quad (30)$$

$$\boldsymbol{\lambda}_{i\alpha}^{k+1} = \boldsymbol{\lambda}_{i\alpha}^k + \rho_{i\alpha}[(1+\epsilon)\boldsymbol{\mu}_i^{k+1} - \mathbf{M}_{i\alpha}\boldsymbol{\mu}_\alpha^{k+1}], \quad (31)$$

where  $i \in V, (i, \alpha) \in E$ .

### 5.4 Complexity and Implementation Details

**Complexity.** In terms of computational complexity, as all other update steps have simple closed-form solutions, the main computational cost lies in updating  $\boldsymbol{\mu}_\alpha$ , which is convex quadratic programming with the probability simplex constraint. Its computational complexity is  $O(|\mathcal{X}_\alpha|^3)$ . As the matrix with the largest size is  $\mathbf{M}_{i\alpha}^\top \mathbf{M}_{i\alpha} \in \mathbb{R}^{|\mathcal{X}_\alpha| \times |\mathcal{X}_\alpha|}$  in LS-LP, the space complexity is  $\mathcal{O}(\sum_{\alpha \in F} (|\mathcal{X}_\alpha|)^2)$ . Both the computational and space complexity of AD3 are similar with LS-LP. More detailed analysis about the computational complexity will be presented in Section 7.5.

**Implementation details.** In each iteration, we use the same value of  $\rho$  for all  $\rho_i$  and  $\rho_{i\alpha}$ . After each iteration, we update  $\rho$  using an incremental rate  $\eta$ , i.e.,  $\rho \leftarrow \rho \times \eta$ . An upper limit  $\rho_{upper}$  of  $\rho$  is also set: if  $\rho$  is larger than  $\rho_{upper}$ , it is not updated anymore. The perturbation  $\epsilon$  is set to  $10^{-5}$ , and  $\rho_{upper}$

<sup>1</sup> <http://sigpromu.org/quadprog/download.php?sid=3wtwk5tb>

can be set as any constant than  $\frac{1}{\epsilon}$ , such as  $2 \times 10^5$ . We utilize two stopping criterion jointly, including: 1) the violation of the local consistency constraint, *i.e.*,  $(\sum_{(i,\alpha) \in E} \frac{\rho_{i\alpha}}{2} \|(1+\epsilon)\boldsymbol{\mu}_i - \mathbf{M}_{i\alpha}\boldsymbol{\mu}_\alpha\|_2^2)^{\frac{1}{2}}$ ; 2) the violation of the equivalence constraint  $(1+\epsilon)\boldsymbol{\mu}_i = \mathbf{v}_i$ , *i.e.*,  $(\sum_{i \in V} \frac{\rho_i}{2} \|(1+\epsilon)\boldsymbol{\mu}_i - \mathbf{v}_i\|_2^2)^{\frac{1}{2}}$ . We set the same threshold  $10^{-5}$  for both criterion. If this two violations are lower than  $10^{-5}$  simultaneously, then the algorithm stops.

## 6 Convergence Analysis

The convergence property of the above ADMM algorithm is demonstrated in Theorem 2. Due to the space limit, the detailed proof will be presented in **Appendix A**.

**Theorem 2** *We suppose that  $\rho$  is set to be larger than a constant, then the variable sequence  $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$  generated by the perturbed ADMM algorithm globally converges to  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ , where  $(\mathbf{y}^*, \mathbf{x}^*)$  is the  $\epsilon$ -KKT point to the LS-LP problem (18), as defined in Definition 2.*

Furthermore, according to Definition 1, we assume that  $\mathcal{L}_{\rho,\epsilon}$  has the KL property at  $(\mathbf{y}^*, \mathbf{x}^*, \boldsymbol{\lambda}^*)$  with the concave function  $\varphi(s) = cs^{1-p}$ , where  $p \in [0, 1)$ ,  $c > 0$ . Consequently, we can obtain the following inequalities:

- (i) If  $p = 0$ , then the perturbed ADMM algorithm will converge in finite steps.
- (ii) If  $p \in (0, \frac{1}{2}]$ , then we will obtain the  $\epsilon$ -KKT solution to the LS-LP problem in at least  $O(\log_{\frac{1}{2}}(\frac{1}{\epsilon})^2)$  steps, with  $\tau \in (0, 1)$  being a small constant, which will be later defined in Appendix A.5.
- (iii) If  $p \in (\frac{1}{2}, 1)$ , then we will obtain the  $\epsilon$ -KKT solution to the LS-LP problem in at least  $O((\frac{1}{\epsilon})^{\frac{4p-2}{1-p}})$  steps.

*Proof* The general structure of the proof consists of the following steps, as follows:

1. The perturbed augmented Lagrangian function  $\mathcal{L}_{\rho,\epsilon}$  (see (24)) is monotonically decreasing along the optimization.
2. The variable sequence  $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$  is bounded.
3. The sequence of variable residuals is converged, *i.e.*,  $\{\|\mathbf{y}^{k+1} - \mathbf{y}^k\|, \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|, \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|\} \rightarrow 0$ , as  $k \rightarrow \infty$ .
4. The variable sequence  $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$  globally converges to the cluster point  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ .
5.  $(\mathbf{y}^*, \mathbf{x}^*)$  is the  $\epsilon$ -KKT point of the LS-LP problem (18).
6. We finally analyze the convergence rate that how many steps are required to achieve the  $\epsilon$ -KKT point.

Table 1: Benchmark datasets used in the Probabilistic Inference Challenge (PIC 2011) [7] and OpenGM 2 [14].  $C_1$  to  $C_7$  represent: number of models, average variables, average factors, average edges, average factor sizes (*i.e.*, the average number of adjacent variables for each factor), average variable states, average factor states.

dataset	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$
Seg-2	50	229.14	622.28	1244.56	2	2	4
Seg-21	50	229.14	622.28	1244.56	2	21	441
Scene	715	182.56	488.99	977.98	2	8	64
Grids	21	3142.86	6236.19	12472.4	2	2	4
Protein	7	14324.7	21854.7	57680.4	2.64	2	6.56

## 7 Experiments

### 7.1 Experimental Settings

#### 7.1.1 Datasets

We evaluate on four benchmark datasets from the Probabilistic Inference Challenge (PIC 2011) [7] and OpenGM 2 [14], including Segmentation [7], Scene [10], Grids [7], and Protein [11], as shown in Table 1. Segmentation consists of Seg-2 and Seg-21, with different variable states. Protein includes higher-order potentials, while others include pairwise potentials.

#### 7.1.2 Compared Methods

We compare with different categories of MAP inference methods, including: **1)** moving making methods, *i.e.*, ICM [3]; **2)** message-passing methods, including belief propagation (BP) [21] and TRBP [38]; **3)** polyhedral methods (including LP relaxation based methods), including dual decomposition using sub-gradient (DD-SG) [15], TRWS [19], AD-Sal [33], PSDD [20] and AD3 [27][28]. **4)** We also compare with LP-LP, which calls the the active-set method (implemented by *linprog* in MATLAB) to optimize  $LP(\theta)$ . It serves as a baseline to measure the performance of above methods. **5)** The most related work  $\ell_2$ -Box ADMM (*i.e.*, the special case of  $\ell_2$ -Box ADMM with  $p = 2$ ) algorithm [41] is also compared. However, the presented algorithm in [41] can only handle MRF models with pairwise potentials, which is formulated as a binary quadratic programming (BQP) problem. Thus,  $\ell_2$ -Box ADMM (hereafter we call it  $\ell_2$ -Box for clarity) is not compared on Protein, of which models include high-order potentials. **6)** We also compare with two hybrid methods, including method DAOOPT (adopting branch-and-bound method [23] as a sub-routine) [31][32] and MPLP-C [37] (adopting MPLP [9] as a sub-routine). The ‘hybrid’ indicates that the method is a combination of an off-the-shelf single method and some heuristic

steps. And we call above 5 types as non-hybrid methods. Both the proposed LS-LP and  $\ell_2$ -Box are implemented by MATLAB. The following methods are implemented by the author provided C++ package, including: PSDD and AD3<sup>2</sup>, MPLP-C<sup>3</sup>, and DAOOPT<sup>4</sup>. All other methods are implemented through the OpenGM 2 software [14], and we add a prefix “ogm” before the method name, such as ogm-TRWS.

In experiments, we set some upper limits: the maximal iteration as 2000 for PSDD and AD3, 500 for  $\ell_2$ -Box and LS-LP, and 1000 for other methods; for DAOOPT, the memory limit of mini buckets is set as 4000 MB and the upper time limit as 2 hours. The parameter tuning of all compared methods (except  $\ell_2$ -Box) is self-included in their implementations. Both LS-LP and  $\ell_2$ -Box are ADMM algorithms, and their hyper-parameters are tuned as follows: the hyper-parameters  $\rho_0$ ,  $\eta$  and  $\rho_{upper}$  (see implementation details of Section 5) are adjusted in the ranges  $\{0.05, 0.1, 1, 5, 10, 10^2, 10^3, 10^4\}$ ,  $\{1.01, 1.03, 1.05, 1.1, 1.2\}$  and  $\{10^6, 10^8\}$ , respectively, and those leading to the higher logPot value are used.

### 7.1.3 Evaluation Metrics

We evaluate the performance of all compared methods using three types of metrics, including the log potential (log-Pot) values, the solution type, as well as the computational complexity and runtime.

**Evaluation using logPot values.** The logPot value indicates the objective value  $\langle \theta, \mu \rangle$  of  $\text{MAP}(\theta)$  (see Eq. (5)). Given that the constraint  $\mathcal{M}_G$  of  $\text{MAP}(\theta)$  is satisfied, the larger logPot value indicates the better inference performance. Since LP-LP gives the optimal solution to  $\text{LP}(\theta)$  (see Eq. (7)) constrained to  $\mathcal{L}_G$ , and we know that  $\mathcal{M}_G \subseteq \mathcal{L}_G$ , then the logPot value of any valid label configurations cannot be larger than that of LP-LP. Note that in the implementation of OpenGM 2 [14], a rounding method is adopted as the post-processing step to produce the integer solution for the continuous MAP inference methods. However, the performance of different MAP inference methods may be significantly changed by rounding. Thus, for other methods not implemented by OpenGM 2, we report the logPot values of original continuous solutions, without any rounding.

**Evaluation using solution types.** Since LP-LP, PSDD and AD3 are possible to give continuous solutions, the larger logPot value doesn’t always mean the better MAP inference result. Thus, we also define four qualitative measures, including *valid*, *uniform*, *fractional* and *approximate*, to intuitively measure the inference quality. *Valid* (V) means that the solution is integer and satisfies the constraints in  $\mathcal{L}_G$ ; *Uniform* (U) denotes that the solution belongs to  $\mathcal{L}_G$ , but

Table 2: A brief illustration of four types of measures for inference quality, on a toy graphical model with two connected variable nodes.

Possible states	Variable 1 {0, 1}	Variable 2 {0, 1}	Factor {00, 01, 10, 11}	Measure
Inferred	(1, 0)	(0, 1)	(0, 1, 0, 0)	Valid
Probability	(0.5, 0.5)	(0.5, 0.5)	(0.25, 0.25, 0.25, 0.25)	Uniform
	(0.2, 0.8)	(0.4, 0.6)	(0.08, 0.12, 0.32, 0.48)	Fractional
	(0.2, 0.8)	(0.4, 0.6)	(0.16, 0.3, 0.4, 0.14)	Approximate

the value is uniform, such as (0.5, 0.5) for the variable node with binary states. *Fractional* (F) indicates that the solution belongs to  $\mathcal{L}_G$ , but its value is fractional, while not uniform; *Approximate* (A) means that some constraints in  $\mathcal{L}_G$  are violated, and its solutions is integer or fractional. Note that it makes sense to compare the logPot values for different solutions, only when the solutions are valid. In contrast, it makes no sense to compare the logPot values, if the solutions belong to the other three types of measures. To better illustrate above measures, we present a brief example on a toy graphical model, as shown in Table 2.

**Evaluation using the computational complexity and practical runtime.** The computational complexity and the practical runtime are also important performance measures for MAP inference methods, as shown in Section 7.5.

## 7.2 Results on Segmentation and Scene

The average results on Seg-2, Seg-21 and Scene are shown in Table 3. LP-LP gives valid solutions on all models, *i.e.*, the best solutions. Except for PSDD, all other methods give valid solutions, and their logPot values can not be higher than those of LP-LP. The logPot values of ICM are the lowest, and those of ogm-BP, ogm-TRBP are slightly lower than the best logPot values, while other methods achieve the best logPot values on most models. Only PSDD gives approximate solutions (*i.e.*, the constraints in  $\mathcal{L}_G$  are not fully satisfied) on some models, specifically, 5 models in Seg-2, 8 models in Seg-21 and 166 models in Scene. ogm-DD-SG fails to give solutions on some models of these datasets, thus we ignore it. Evaluations on these easy models only show that the performance ranking is ogm-ICM < ogm-BP, ogm-TRBP,  $\ell_2$ -Box < others.

## 7.3 Results on Grids

The results on Grids are shown in Table 4. For clarity, we use the model indexes M1 to M21 to indicate the model name to save space in this section. The corresponding model names from M1 to M21 are grid20x20.f10.uai, grid20x20.f10.wrap.uai, grid20x20.f15.uai, grid20x20.f15.wrap.uai, grid20x20.f5.wrap.uai, grid40x40.f10.uai, grid40x40.f10.wrap.uai, grid

<sup>2</sup> <http://www.cs.cmu.edu/ark/AD3/>

<sup>3</sup> <https://github.com/openGM/MPLP>

<sup>4</sup> <https://github.com/lotten/daoopt>



Table 3: LogPot values of MAP inference solutions on Seg-2, Seg-21 and Scene. Except of PSDD, all other methods give valid solutions. The best result among valid solutions in each row is highlighted in bold. Please refer to Section 7.2 for details.

Method type → Dataset ↓		Baseline LP-LP	Hybrid methods		Non-hybrid methods							Proposed LS-LP	
			DAOOPT	MPLP-C	ogm-ICM	ogm-BP	ogm-TRBP	ogm-TRWS	ogm-ADSal	PSDD	AD3	$\ell_2$ -Box	
Seg-2	mean	<b>-75.5</b>	<b>-75.5</b>	<b>-75.5</b>	-137.1	-79	-76.8	<b>-75.5</b>	<b>-75.5</b>	-75.4	<b>-75.5</b>	-76.5	-75.6
	std	19.63	19.63	19.63	70.1	20.24	19.36	19.24	19.24	19.77	19.63	20.3	19.69
Seg-21	mean	<b>-324.89</b>	-325.34	<b>-324.89</b>	-393.37	-330.37	-328.92	<b>-324.89</b>	<b>-324.89</b>	-325.1	<b>-324.89</b>	-344.51	<b>-324.89</b>
	std	58.12	58.14	58.12	74.47	58.54	58.57	56.97	56.97	58.16	58.12	59.24	58.12
Scene	mean	<b>866.66</b>	<b>866.66</b>	<b>866.66</b>	864.27	866.49	866.51	<b>866.66</b>	<b>866.66</b>	866.65	<b>866.66</b>	864.11	<b>866.66</b>
	std	109.34	109.34	109.36	109.64	109.22	109.2	109.19	109.19	109.34	109.34	108.66	109.34

Table 4: MAP inference results on Grids dataset. LP-LP, PSDD and AD3 produce uniform solutions on all models in Grids, while all other methods give valid solutions. Here we only show the logPot of LP-LP as the upper bound of other methods. The best logPot among integer solutions in each row is highlighted in bold. The number with in circle indicates the performance ranking of each method. Please refer to Section 7.3 for details.

Method type → Model ↓		Baseline LP-LP	Hybrid methods		Non-hybrid methods							Proposed LS-LP
			DAOOPT	MPLP-C	ogm-ICM	ogm-BP	ogm-TRBP	ogm-DD-SG	ogm-TRWS	ogm-ADSal	$\ell_2$ -Box	
M1	3736.7	<b>3015.7</b> ①	<b>3015.7</b> ①	2708.9 ⑤	121.3 ⑥	-235.2 ⑩	1286.3 ⑥	2524.9 ⑦	2605.2 ⑥	2794.8 ④	2931.8 ③	
M2	3830.3	<b>3051</b> ①	3033.6 ②	2567.9 ⑦	276.4 ⑥	19.2 ⑩	1484.7 ⑥	2674.4 ⑥	2670.2 ⑥	2812.4 ④	2936.7 ③	
M3	5605.1	<b>4517.3</b> ①	<b>4517.3</b> ①	4067.3 ③	332.1 ⑥	14.02 ⑩	1889.7 ⑥	3829.3 ⑦	3884 ⑥	4301.1 ④	4408.9 ③	
M4	5745.5	<b>4563.2</b> ①	<b>4563.2</b> ①	3837.12 ②	924.5 ⑥	-36.7 ⑩	2023.4 ⑥	3894.6 ⑥	4015 ⑥	4202.4 ④	4446.6 ③	
M5	1915.2	<b>1542.7</b> ①	<b>1542.7</b> ①	1318.41 ⑦	481.5 ⑥	-47.8 ⑩	807.6 ⑥	1325.5 ⑤	1323.9 ⑥	1427.4 ④	1503.2 ③	
M6	15601.2	12662.9 ⑥	<b>12665.7</b> ①	10753.7 ⑥	2793.5 ⑥	2214.3 ⑩	5051.9 ⑥	10500.8 ⑦	11029 ⑤	11486.2 ⑥	12336.1 ⑤	
M7	16291.5	13050.7 ⑥	<b>13054.8</b> ①	10903.8 ⑤	1217.1 ⑥	132.4 ⑩	4634.8 ⑥	10665 ⑦	10870.4 ⑥	11867.6 ⑥	12537.2 ⑤	
M8	23401.8	<b>18952.45</b> ①	18896.8 ②	16154.2 ②	4314.9 ⑩	5371.1 ⑥	7160 ⑥	16014 ⑦	16276.9 ⑤	17367.5 ④	18358.7 ③	
M9	24437.3	<b>19538</b> ①	19427.5 ②	16334.2 ②	3560.8 ⑥	-1111 ⑩	7187.3 ⑥	16004.3 ⑦	16508.1 ⑤	17990 ④	18785.8 ③	
M10	3121.2	<b>2689</b> ①	2688.8 ②	2255.38 ⑥	1665.3 ⑥	1582.9 ⑥	1330.7 ⑩	2215.7 ⑦	2369.1 ⑤	2552.6 ④	2659.8 ③	
M11	3231.6	<b>2714.67</b> ①	2714.52 ②	2258.54 ⑦	1399.6 ⑥	42.8 ⑩	1285.9 ⑥	2271.5 ⑥	2370.1 ⑤	2556.8 ④	2654.9 ③	
M12	7800.6	<b>6401.15</b> ①	6396 ②	5356.28 ⑥	2033.5 ⑥	1953.1 ⑩	2832.5 ⑥	5282.5 ⑦	5558.8 ⑥	5903 ④	6201.2 ③	
M13	8078.5	<b>6472.9</b> ①	6469.7 ②	5425.16 ⑦	1711.5 ⑥	381 ⑩	2814.3 ⑥	5452.8 ⑥	5646.1 ⑤	5923.5 ④	6275.4 ③	
M14	62943	–	45813.6 ②	43538.9 ②	5690.9 ⑥	6426.7 ⑥	18700.4 ⑦	42274.2 ⑥	43292.5 ⑤	44397.5 ④	<b>48766.1</b> ①	
M15	63993.1	–	47444.4 ②	42855 ③	4287.1 ⑥	956.4 ⑥	18811.9 ⑦	42535 ⑥	42918.7 ④	44759.5 ④	<b>48657.3</b> ①	
M16	94414.5	–	69408.6 ②	65081.2 ②	4374.2 ②	4656.5 ⑥	27320.6 ⑦	63148.1 ⑥	64401.1 ⑤	66784.2 ④	<b>72993.8</b> ①	
M17	96243.6	–	71730.8 ②	63768.1 ②	13662.7 ⑥	-529.3 ⑥	27287.7 ⑦	63885.1 ⑤	64487.9 ④	67589.4 ④	<b>73486</b> ①	
M18	12721.3	–	10445.8 ②	9062.03 ③	5198.7 ⑥	4975.4 ⑥	4785.5 ⑥	8793.5 ⑥	9408.4 ④	10015.1 ④	<b>10580.8</b> ①	
M19	12875.6	–	10674.1 ②	9214.57 ③	5944.6 ⑥	1213.1 ⑥	5328.5 ⑥	8952.4 ⑥	9385.1 ④	10163.6 ④	<b>10698.4</b> ①	
M20	31809.7	–	22292.5 ③	21527.9 ③	5410.9 ⑥	4762 ⑥	9837.3 ⑦	21546.8 ⑤	22109.5 ④	22913.3 ④	<b>24834.5</b> ①	
M21	31996.9	–	24032.4 ③	21529.6 ③	4242.3 ③	47.6 ⑥	10423.8 ⑦	21195.9 ⑥	21730.3 ④	22668.7 ④	<b>24532.8</b> ①	

40x40.f15.uai, grid40x40.f15.wrap.uai, grid40x40.f2.uai, grid 40x40.f2.wrap.uai, grid40x40.f5.uai, grid40x40.f5.wrap.uai, grid80x80.f10.uai, grid80x80.f10.wrap.uai, grid80x80.f15.uai, grid 80x80.f15.wrap.uai, grid80x80.f2.uai, grid80x80.f2.wrap.uai, grid80x80.f5.uai, grid80x80.f5.wrap.uai, respectively.

The models in Grids are much challenging for LP relaxation based methods, as all models have symmetric pairwise log potentials and very dense cycles in the graph. In this case, many vertices of  $\mathcal{L}_G$  are uniform solutions (0.5, 0.5). Consequently, the LP relaxation based methods are likely to produce uniform solutions. This is verified by that LP-LP, AD3, PSDD give uniform solutions on all models in Grids, *i.e.*, most solutions are 0.5. Thus, we only show the logPot values of LP-LP in Table 4, to provide the theoretical upper-bound of logPot of valid solutions from other methods. In contrast, the additional  $\ell_2$ -sphere constraint in LS-LP excludes the uniform solutions. On small scale models M1

to M13, DAOOPT and MPLP-C show the highest logPot values, while LS-LP gives slightly lower values. On large scale models M14 to M21, DAOOPT fails to give any result within 2 hours. LS-LP gives the best results, while MPLP-C shows slightly lower results.  $\ell_2$ -Box performs worse than LS-LP, MPLP-C and DAOOPT on most models, while better than all other methods, among which ogm-BP, ogm-TRBP and ogm-DD-SG perform worst. These results demonstrate that **1)** LS-LP is comparable to hybrid methods DAOOPT and MPLP-C, but with much lower computational cost (shown in Section 7.5); **2)** LS-LP performs much better than other non-hybrid methods.

Table 5: LogPot values of MAP inference solutions on Protein dataset. Except for PSDD and AD3, all other methods give integer solutions. Both PSDD and AD3 produce mixed types of solutions on all models. The best result among valid solutions in each row is highlighted in bold. The number with in circle indicates the performance ranking of each method. Please refer to Section 7.4 for details.

Method type $\rightarrow$	Hybrid methods	Non-hybrid methods						Proposed
Model $\downarrow$	MPLP-C	ogm-ICM	ogm-BP	ogm-TRBP	ogm-DD-SG	PSDD	AD3	LS-LP
M1	-30181.3 ②	-32409.9 ⑤	-32019.1 ④	-31671.6 ③	-33381.2 ⑥	-30128.8	-30143.6	<b>-30165.5 ①</b>
M2	-29305.4 ②	-32561.3 ⑤	-30966.1 ③	-31253.3 ④	-33583.6 ⑥	-29307.3	-29302.6	<b>-29295.4 ①</b>
M4	-28952.1 ②	-32570 ⑤	-31031.4 ③	-31176.6 ④	-33747.7 ⑥	-28952.5	-28952	<b>-28952 ①</b>
M5	-269567 ③	<b>-256489 ①</b>	-382766 ⑤	-357330 ④	-553376 ⑥	-66132.3	-115562	-267814 ②
M6	-30070.6 ②	-31699.1 ⑤	-30765.2 ③	-30772.2 ④	-32952.9 ⑥	-30063.6	-30062.2	<b>-30063.4 ①</b>
M7	-30288.3 ②	-32562.2 ⑤	-31659.6 ③	-31791.1 ④	-33620.4 ⑥	-30248.5	-30239.8	<b>-30266 ①</b>
M8	-29336.5 ②	-32617.2 ⑤	-31064.7 ③	-31219.9 ④	-34549.9 ⑥	-29331	-29336.1	<b>-29334.7 ①</b>

## 7.4 Results on Protein

The results on Protein are shown in Table 5. Different with above three datasets, Protein includes 8 large scale models, and with high-order factors. Similarly, we use the model indexes M1 to M8 to indicate the model name to save space in this section. The corresponding model names of M1 to M8 are didNotconverge1.uai, didNotconverge2.uai, didNotconverge4.uai, didNotconverge5.uai, didNotconverge6.uai, didNotconverge7.uai, didNotconverge8.uai, respectively. As M1 and M3 are the same model, we remove M3 in experiments. DAOOPT fails to give solutions within 2 hours on all models, and LP-LP cannot produce solutions due to the memory limit. ogm-TRWS and  $\ell_2$ -Box are not evaluated as it cannot handle high-order factors. LS-LP produces valid integer solutions on all models, and gives the highest logPot values on all models except of M5. MPLP-C gives slightly lower logPot values than LS-LP. AD3 only produces a fractional solution on M4, while produces approximate and fractional solutions on other models, while PSDD gives approximate and fractional solutions on all models. Specifically, the solution types of AD3 on M1 to M8 are:  $A + 2.01\%U + 21.52\%F$ ;  $A + 0.55\%U + 0.66\%F$ ;  $0.17\%U$ ;  $A + 9.51\%U + 32.44\%F$ ;  $A + 0.74\%U + 13.7\%F$ ;  $A + 3.41\%U + 17.16\%F$ ;  $A + 0.38\%U + 0.13\%F$ . Those of PSDD are:  $A + 6.96\%U + 11.75\%F$ ;  $A + 1.47\%U + 3.19\%F$ ;  $A + 0.73\%U + 2.06\%F$ ;  $A + 13.72\%U + 19.67\%F$ ;  $A + 4.35\%U + 7.9\%F$ ;  $A + 6.28\%U + 10.4\%F$ ;  $A + 0.68\%U + 1.89\%F$ . Other methods also show much worse performance than LS-LP and MPLP-C. One exception is that ogm-ICM gives the best results on M5, and we find that M5 is the most challenging model for approximated methods.

Table 6: Computational complexities of all compared methods. Excluding  $\mathcal{E}$ , the definitions of all other notations can be found in Section 3.  $\mathcal{E}$  denotes the edge set of the original MRF graph, while  $E$  indicates the edge set of the corresponding factor graph.  $T$  represents the number of iterations.

Methods	Complexities
MPLP	$O(\sum_i^V  \mathcal{N}_i ^2 \cdot  \mathcal{X}_i  + 2 \sum_{(i,j) \in \mathcal{E}} ( \mathcal{X}_i  +  \mathcal{X}_j ))$
MPLP-C	$O(T_{\text{outer}} [T_{\text{inner}} O(\text{MPLP}) +  \mathcal{E} ])$
ogm-ICM	$O(T [\sum_i^V  \mathcal{X}_i ])$
ogm-BP	$O(T [\sum_i^V ( \mathcal{N}_i  - 1) \sum_{\alpha}^{\mathcal{N}_i}  \mathcal{X}_{\alpha}  + \sum_{\alpha}^F ( \mathcal{N}_{\alpha}  - 1) \sum_i^{\mathcal{N}_{\alpha}}  \mathcal{X}_i ])$
ogm-TRWS	$O(T [ \mathcal{E}  \cdot \max_{i \in V}  \mathcal{X}_i ])$
ogm-ADSal	
PSDD	$O(T [\sum_i^V ( \mathcal{N}_i  \cdot  \mathcal{X}_i  + \sum_{\alpha}^F  \mathcal{X}_{\alpha} ^3 + \sum_{(i,\alpha) \in E}  \mathcal{X}_i  \cdot  \mathcal{X}_{\alpha} ])$
AD3	
$\ell_2$ -Box	$O(T [\sum_i^V  \mathcal{X}_i ^3])$
LS-LP	$O(T [\sum_i^V  \mathcal{X}_i  + \sum_{\alpha}^F  \mathcal{X}_{\alpha} ^3 + \sum_{(i,\alpha) \in E}  \mathcal{X}_i  \cdot  \mathcal{X}_{\alpha} ])$

## 7.5 Comparisons on Computational Complexities and Practical Runtime

**Computational complexities** of all compared methods (except of LP-LP and DAOOPT) are summarized in Table 6. As there is no clear conclusion of the complexity of the active-set algorithm for linear programming, the complexity of LP-LP is not presented. DAOOPT [32] is combination of 6 sequential sub-algorithms and heuristic steps, thus its computational complexity cannot be computed. As these complexities depend on the graph structure (*i.e.*,  $V, E, \mathcal{E}, F, \mathcal{N}_i, \mathcal{N}_{\alpha}, \mathcal{X}_i, \mathcal{X}_{\alpha}$ ), it is impossible to give a fixed ranking of them. However, it is notable that the complexity of LS-LP is linear w.r.t. the number of variables, factors and edges (*i.e.*,  $|V|, |F|, |E|$ ). Moreover, as the sub-problems w.r.t. variables/factors/edges in LS-LP are independent, they can be solved in parallel. However, the complexity of LS-LP is super-linear w.r.t. the state space  $|\mathbf{X}_{\alpha}|$  and  $|\mathbf{X}_i|$ . Thus, the proposed LS-LP method is suitable for graphical models with very large-

Table 7: Iterations and practical runtime on Seg-2, Seg-21 and Scene.

Datasets		LP-LP	DAOOPT	MPLP-C	ogm-ICM	ogm-BP	ogm-TRBP	ogm-TRWS	ogm-ADSal	PSDD	AD3	$\ell_2$ -Box	LS-LP	
Seg-2	iters	mean	9.4	–	1.3	2.12	1000	1000	18.74	16.96	632.5	70.12	78.84	84.14
		std	1.03	–	0.463	0.32	0	0	15.67	8.55	618.2	76.79	100.6	49.58
	runtime	mean	0.096	0.54	0.033	0.007	36.51	40.73	0.191	0.7828	0.032	0.001	0.849	12.53
		std	0.01	0.504	0.054	0.001	0.61	0.67	0.167	0.6373	0.027	0.006	1.044	7.39
Seg-21	iters	mean	16.13	–	1.46	76.9	1000	1000	23.4	17.58	783	73.26	39.4	78.63
		std	1.55	–	0.646	82.2	0	0	20.53	11.95	672	45.09	68.1	34.65
	runtime	mean	45.71	1311.8	0.745	0.185	1612	1891.6	16.33	57.01	1.14	0.228	8.643	342
		std	9.34	1593.2	1.1	0.062	25.05	25.7	12.89	44.51	1.1	0.114	14.71	147
Scene	iters	mean	12.56	–	1.32	319.5	1000	1000	15.03	10.94	791	71.59	43.03	75.36
		std	1.04	–	0.48	58.89	0	0	12.21	11.15	766	82.83	37.47	52.47
	runtime	mean	1.907	82.89	0.072	0.081	229.67	265.45	1.7	5.03	0.066	0.029	0.777	29.23
		std	0.266	18.46	0.111	0.011	15.46	17.6	1.3	5.88	0.071	0.024	0.681	20.36

Table 8: Iterations and practical runtime on Grids.

Models		DAOOPT	MPLP-C	ogm-ICM	ogm-BP	ogm-TRBP	ogm-DD-SG	ogm-TRWS	ogm-ADSal	$\ell_2$ -Box	LS-LP
M1	iters	–	420	195	1000	1000	1000	14	39	22	234
	runtime	8	303.9	0.012	40.1	40.89	86.9	0.13	2.79	0.317	41.4
M2	iters	–	1000	192	1000	1000	1000	18	37	35	227
	runtime	10	613.3	0.013	42.71	42.13	91.1	0.18	2.72	0.478	39.4
M3	iters	–	686	195	1000	1000	1000	13	36	36	155
	runtime	8	331.5	0.022	40.2	39.28	86.4	0.12	2.52	0.412	27
M4	iters	–	1000	193	1000	1000	1000	13	42	27	434
	runtime	10	523.5	0.013	42.72	45.53	91.5	0.17	3.16	0.318	75.3
M5	iters	–	1000	193	1000	1000	1000	21	37	26	265
	runtime	11	468.7	0.013	42.77	48.07	91.2	0.02	2.73	0.291	45.7
M6	iters	–	1000	741	1000	1000	1000	19	37	18	271
	runtime	7200	840.3	0.052	166.6	182.3	362.3	0.78	11.1	2.19	46.5
M7	iters	–	1000	750	1000	1000	1000	18	37	16	204
	runtime	7200	1040.5	0.053	171.6	187.2	371.7	0.7	11.4	1.96	35.3
M8	iters	–	1000	730	1000	1000	1000	17	37	29	407
	runtime	7200	917.1	0.072	166.6	182.1	362.4	0.68	10.9	3.47	66.1
M9	iters	–	1000	739	1000	1000	1000	16	36	31	327
	runtime	7200	1096.8	0.053	171.7	189.5	372.3	0.85	10.6	3.80	56.3
M10	iters	–	521	738	1000	1000	1000	45	38	45	500
	runtime	7200	273.9	0.072	166.6	182.1	362.4	2.14	10.6	5.68	88.3
M11	iters	–	1000	765	1000	1000	1000	36	44	39	390
	runtime	7200	843.1	0.053	171.7	187.7	372.2	1.57	12.6	5.13	68.4
M12	iters	–	1000	727	1000	1000	1000	26	38	19	314
	runtime	7200	824.8	0.073	166.7	187.6	361.5	1.05	11	2.25	52.5
M13	iters	–	1000	755	1000	1000	1000	34	41	17	493
	runtime	7200	999.7	0.052	171.8	188.4	372.3	1.44	11.9	2.14	85.9
M14	iters	–	1000	3029	1000	1000	1000	21	37	12	290
	runtime	–	1956.4	0.218	676.5	756	1474	3.52	43.5	21.72	50.5
M15	iters	–	1000	2991	1000	1000	1000	19	36	12	202
	runtime	–	2028.9	0.218	687.2	770.2	1495	3.37	42	21.65	33.7
M16	iters	–	1000	3038	1000	1000	1000	15	35	20	435
	runtime	–	2110.4	0.217	676.8	762.2	1471	2.58	41.5	36.53	76
M17	iters	–	1000	3048	1000	1000	1000	20	37	15	354
	runtime	–	2238.8	0.221	686.8	767.6	1492.2	3.37	43	26.93	64
M18	iters	–	1000	3173	1000	1000	1000	54	53	35	250
	runtime	–	1956.3	0.221	677	756.1	1475.2	8.87	63.7	60.28	43.5
M19	iters	–	1000	3095	1000	1000	1000	109	45	37	250
	runtime	–	1953.5	0.221	687.4	773	1493.8	19.13	50.7	63.93	39.9
M20	iters	–	1000	3038	1000	1000	1000	30	38	13	316
	runtime	–	1974.2	0.206	676.4	756	1473.2	4.99	42.5	23.06	55.4
M21	iters	–	1000	2985	1000	1000	1000	22	37	12	280
	runtime	–	2167.3	0.218	687.2	767.5	1494.4	3.75	40.6	21.54	49.6

Table 9: Iterations and practical runtime on Protein.

Models		MPLP-C	ogm-ICM	ogm-BP	ogm-TRBP	ogm-DD-SG	PSDD	AD3	LS-LP
M1	iters	499	797	1000	1000	1000	2000	2000	161
	runtime	2320	0.39	32019	31671	33381	14.14	20.28	1746
M2	iters	478	634	1000	1000	1000	2000	2000	98
	runtime	2399	0.44	30966	31253	33583	7.74	3.87	1076
M4	iters	24	859	1000	1000	1000	2000	1524	140
	runtime	80	0.51	31031	31177	33748	5.73	3.32	1535
M5	iters	500	5275	1000	1000	1000	2000	2000	1000
	runtime	2248	0.8	382766	357330	553376	22.62	28.36	10787
M6	iters	433	597	1000	1000	1000	2000	2000	482
	runtime	2427	0.36	30765	30772	32953	12.96	14.12	5314
M7	iters	475	836	1000	1000	1000	2000	2000	145
	runtime	2399	0.43	31660	31791	33620	13.55	18.63	1989.7
M8	iters	16	778	1000	1000	1000	2000	2000	186
	runtime	80	0.51	31065	31220	34550	5.13	3.31	1566

scale variables and dense connections, but with modest state spaces for variables and factors.

**Practical runtime.** Due to the dependency of the computational complexity on the graph structure, the practical runtime of these methods will vary significantly on different graphs. In the following, we present the practical runtime of all compared methods on above evaluated datasets. To test the runtime fairly, we run all methods at the same machine, and only run one experiment at the same time. The iterations and practical runtime on different datasets are shown respectively in Table 7 for Segmentation and Scene, Table 8 for Grids and Table 9 for Protein.

As shown in Table 7, on the small and easy models, the runtime of LP-LP, MPLP-C, ogm-ICM, PSDD and AD3 are very small, and the runtime of ogm-BP, ogm-TRBP, ogm-TRWS, ogm-ADSaI, LS-LP and  $\ell_2$ -Box are larger, while the runtime of DAOOPT are the largest. Besides, the iterations of AD3 and LS-LP are much smaller than the one of PSDD, given the fact that their similar computational complexities per iteration. The iterations of LP-LP and MPLP-C are also provided, but they are incomparable with PSDD, AD3 and LS-LP. The complexity of each iteration in LP-LP depends on the problem size  $|\mu|$  (see Eq. (7)). In MPLP-C, each outer iteration includes 100 iterations of MPLP and adding violated constraints. The complexity of each MPLP iteration is stable, but the complexity of adding violated constraints varies significantly in different outer iterations.

In terms of the comparison on Grids (see Table 8), the iterations and runtime of LP-LP are smaller than those of other methods, but it gives uniform solutions on all models. As demonstrated Section 7.2, LP-LP, PSDD and AD3 produce uniform solutions on all models, thus we also don't present their iterations and runtime. The runtime of DAOOPT on small models (*i.e.*, M1 to M5) are small, and it achieves 7200 seconds (the upper limit) on M6 to M13, while it cannot give any solution in 7200 seconds on M14 to M21. For

MPLP-C, both iterations and runtime are large on all models. Note that the runtime per iteration of MPLP-C becomes larger along with the model scale. For ICM, the iteration is large on multiple models, but with very small runtime, as its complexity per iteration is low. Both message passing methods, including ogm-BP and ogm-TRBP, achieve the upper limit of iterations. It demonstrate that their convergence is very slow. The convergence of ogm-DD-SG is also slow, and its runtime per iteration is even higher than above two message passing methods. Two LP relaxation based methods, including ogm-TRWS and ogm-ADSaI, converge in a few iterations, and are of small runtime. In contrast, the iterations of LS-LP are only larger than those of ogm-TRWS, ogm-ADSaI and  $\ell_2$ -Box, while smaller than other methods. And, the runtime per iteration of LS-LP is similar with that of ogm-BP, while smaller than those of other methods except of ogm-ICM and ogm-TRWS. Considering that C++ is much more efficient than MATLAB, if LS-LP is also implemented by C++, its runtime should be much lower than those of most compared methods. In other words, the computational complexity of LS-LP is much smaller than most compared methods. Note that the runtime per iteration of  $\ell_2$ -Box increases along with the model size. For example, its runtime per iteration on M1 is 0.014 seconds, while that on M21 is 1.795 seconds. In contrast, the runtime per iteration of LS-LP are rather stable. As shown in Table 6, the complexity per iteration of  $\ell_2$ -Box is  $O([\sum_i^V |\mathcal{X}_i|]^3)$ . Obviously,  $\ell_2$ -Box is difficult to apply to large-scale models. In contrast, LS-LP is conducted based on the decomposition of the factor graph to independent factors and variables, due to which the parallel computation is allowed. Thus, the scalability of LS-LP is much better than  $\ell_2$ -Box for MAP inference.

In terms of the comparison on Protein (see Table 9), the ascending ranking of runtime is ogm-ICM, PSDD, AD3, MPLP-C, LS-LP, ogm-BP, ogm-TRBP, ogm-DD-SG. Although

the runtime of PSDD and AD3 are very small, but they only give *approximate* solutions in 2000 iterations on all models, except for AD3 on M4. For MPLP-C, the iterations and runtime are small on M4 and M8, but very large on all other 5 models. In contrast, although the runtime of LS-LP is much larger than those of PSDD and AD3, its iterations are much smaller on all models. If LS-LP is also implemented by C++, its practical runtime will be much smaller than those of PSDD and AD3.

In summary, the above comparisons on iterations and practical runtime demonstrate:

1. LS-LP converges much faster than most compared methods, except of ogm-TRWS, ogm-ADSal and  $\ell_2$ -Box. On large-scale models (see Table 5), the complexities per iteration of LS-LP, PSDD and AD3 are similar, and are much smaller than those of other methods (except of ICM).
2. The complexity of MPLP-C is larger than PSDD, AD3 and LS-LP, while smaller than DAOOPT. But its practical iterations and runtime vary significantly on different models. Both the complexity and practical runtime of DAOOPT are much larger than other methods.
3. Considering the performance of the MAP inference results presented in Section 7.2, 7.3 and 7.4, we conclude that LS-LP shows very competitive performance compared to state-of-the-art MAP inference methods.

## 7.6 Discussions

We obtain three conclusions from above experiments evaluated on different types of models. **1)** Compared with the hybrid methods including DAOOPT and MPLP-C, the performance of LS-LP is comparable. However, the computational cost of LS-LP is much lower, as the hybrid methods adopt LP relaxation based methods as sub-routines. **2)** Compared with LP relaxation based methods, especially PSDD and AD3, LS-LP always give valid solutions, without rounding; and, the logPot values of LS-LP are much higher, with similar computational cost. **3)** Compared to  $\ell_2$ -Box, which can be only applied to models with pairwise potentials, LS-LP is applied to any type of models. Besides, the decomposition of the factor graph allows for the parallel computations with respect to each factor and each variable, while  $\ell_2$ -Box solves a QP problem over the whole MRF model. Thus, LS-LP is a much better choice than  $\ell_2$ -Box for MAP inference. **4)** Compared to other approximated methods, LS-LP always shows much better performance in difficult models (*e.g.*, Grids and Protein).

## 8 Conclusions

In this work, we proposed an novel formulation of MAP inference, called  $\ell_2$ -sphere linear program (LS-LP). Start-

ing from the standard linear programming (LP) relaxation, we added the  $\ell_2$ -sphere constraint onto variable nodes. The intersection between the  $\ell_2$ -sphere constraint and the local marginal polytope  $\mathcal{L}_G$  in LP relaxation is proved to be the exact set of all valid integer label configurations. Thus, the proposed LS-LP problem is equivalent to the original MAP inference problem. By adding a sufficiently small perturbation  $\epsilon$  onto the objective function and constraints, we proposed a perturbed ADMM algorithm for optimizing the LS-LP problem. Although the  $\ell_2$ -sphere constraint is non-convex, we proved that the ADMM algorithm will globally converge to the  $\epsilon$ -KKT point of the LS-LP problem. The analysis of convergence rate is also presented. Experiments on three benchmark datasets show the competitive performance of LS-LP compared to state-of-the-art MAP inference methods.

## A Convergence Analysis

To facilitate the convergence analysis, here we rewrite some equations and notations firstly defined in Sections 5. Problem (11) can be simplified to the following general shape, as follows

$$\text{LS-LP}(\theta) = \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + h(\mathbf{y}), \text{ s.t. } \mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{y}. \quad (32)$$

Our illustration for (18) is separated into three parts, as follows:

1. **Variables.**  $\mathbf{x} = [\boldsymbol{\mu}_1; \dots; \boldsymbol{\mu}_{|V|}] \in \mathbb{R}^{\sum_i |\mathcal{X}_i|}$ , and it concatenates all variable nodes  $\boldsymbol{\mu}_V$ .  $\mathbf{y} = [\mathbf{y}_1; \dots; \mathbf{y}_{|V|}]$  with  $\mathbf{y}_i = [\mathbf{v}_i; \boldsymbol{\mu}_{\alpha_{i,1}}; \dots; \boldsymbol{\mu}_{\alpha_{i,|\mathcal{N}_i|}}] \in \mathbb{R}^{|\mathcal{X}_i| + \sum_{\alpha \in \mathcal{F}} |\mathcal{X}_\alpha|}$ .  $\mathbf{y}$  concatenates all factor nodes  $\boldsymbol{\mu}_V$  and the extra variable nodes  $\mathbf{v}$ ;  $\mathbf{y}_i$  concatenates the factor nodes and the extra variable node connected to the  $i$ -th variable node  $\boldsymbol{\mu}_i$ .  $\mathcal{N}_i$  indicates the set of neighborhood factor nodes connected to the  $i$ -th variable node; the subscript  $\alpha_{i,j}$  indicates the  $j$ -th factor connected to the  $i$ -th variable, with  $i \in V$  and  $j \in \mathcal{N}_i$ .
2. **Objective functions.**  $f(\mathbf{x}) = \mathbf{w}_x^\top \mathbf{x}$  with  $\mathbf{w}_x = -[\boldsymbol{\theta}_1; \dots; \boldsymbol{\theta}_{|V|}]$ ,  $h(\mathbf{y}) = g(\mathbf{y}) + \mathbf{w}_y^\top \mathbf{y}$ , with  $\mathbf{w}_y = [\mathbf{w}_1; \dots; \mathbf{w}_{|V|}]$  with  $\mathbf{w}_i = -[\mathbf{0}; \frac{1}{|\mathcal{N}_{\alpha_{i,1}}|} \boldsymbol{\theta}_{\alpha_{i,1}}; \dots; \frac{1}{|\mathcal{N}_{\alpha_{i,|\mathcal{N}_i|}}|} \boldsymbol{\theta}_{\alpha_{i,|\mathcal{N}_i|}}]$ , and  $\mathcal{N}_\alpha = \{i \mid (i, \alpha) \in E\}$  being the set of neighborhood variable nodes connected to the  $\alpha$ -th factor.  $g(\mathbf{y}) = \mathbb{I}(\mathbf{v} \in \mathcal{S}) + \sum_{\alpha \in \mathcal{F}} \mathbb{I}(\boldsymbol{\mu}_\alpha \in \Delta^{|\mathcal{X}_\alpha|})$ , with  $\mathbb{I}(a)$  being the indicator function:  $\mathbb{I}(a) = 0$  if  $a$  is true, otherwise  $\mathbb{I}(a) = \infty$ .
3. **Constraint matrices.** The constraint matrix  $\mathbf{A} = \text{diag}(\mathbf{A}_1, \dots, \mathbf{A}_i, \dots, \mathbf{A}_{|V|})$  with  $\mathbf{A}_i = [\mathbf{I}_{|\mathcal{X}_i|}; \dots; \mathbf{I}_{|\mathcal{X}_i|}] \in \{0, 1\}^{(|\mathcal{N}_i|+1)|\mathcal{X}_i| \times |\mathcal{X}_i|}$ .  $\mathbf{B} = \text{diag}(\mathbf{B}_1, \dots, \mathbf{B}_i, \dots, \mathbf{B}_{|V|})$ , with  $\mathbf{B}_i = \text{diag}(\mathbf{I}_{|\mathcal{X}_i|}, \mathbf{M}_{i,\alpha_{i,1}}, \dots, \mathbf{M}_{i,\alpha_{i,|\mathcal{N}_i|}})$ .  $\mathbf{A}$  summarizes all constraints on  $\boldsymbol{\mu}_V$ , while  $\mathbf{B}$  collects all constraints on  $\boldsymbol{\mu}_F$  and  $\mathbf{v}$ .

Note that Problem (18) has a clear structure with two groups of variables, corresponding the augmented factor graph (see Fig. 1(c)).

According to the analysis presented in [40], a sufficient condition to ensure the global convergence of the ADMM algorithm for the problem LS-LP( $\theta$ ) is that  $\text{Im}(\mathbf{B}) \subseteq \text{Im}(\mathbf{A})$ , with  $\text{Im}(\mathbf{A})$  being the image of  $\mathbf{A}$ , *i.e.*, the column space of  $\mathbf{A}$ . However,  $\mathbf{A}$  in (32) is full column rank, rather than full row rank, while  $\mathbf{B}$  is full row rank. To satisfy this sufficient condition, we introduce a sufficiently small perturbation to both the objective function and the constraint in (32), as follows

$$\text{LS-LP}(\theta; \epsilon) = \min_{\hat{\mathbf{x}}, \mathbf{y}} \hat{f}(\hat{\mathbf{x}}) + h(\mathbf{y}), \text{ s.t. } \hat{\mathbf{A}}\hat{\mathbf{x}} = \mathbf{B}\mathbf{y}, \quad (33)$$

where  $\hat{\mathbf{A}} = [\mathbf{A}, \epsilon \mathbf{I}]$  with a sufficiently small constant  $\epsilon > 0$ , then  $\hat{\mathbf{A}}$  is full row rank.  $\hat{\mathbf{x}} = [\mathbf{x}; \bar{\mathbf{x}}]$ , with  $\bar{\mathbf{x}} = [\bar{\mathbf{x}}_1; \dots; \bar{\mathbf{x}}_{|V|}] \in \mathbb{R}^{\sum_i (|\mathcal{N}_i|+1)|\mathcal{X}_i|}$  and  $\bar{\mathbf{x}}_i = [\boldsymbol{\mu}_i; \dots; \boldsymbol{\mu}_i] \in \mathbb{R}^{(|\mathcal{N}_i|+1)|\mathcal{X}_i|}$ .  $\hat{f}(\hat{\mathbf{x}}) = f(\mathbf{x}) + \frac{1}{2} \epsilon \bar{\mathbf{x}}^\top \bar{\mathbf{x}}$ . Consequently,  $\text{Im}(\hat{\mathbf{A}}) \equiv \text{Im}(\mathbf{B}) \subseteq \mathbb{R}^{\text{rank of } \hat{\mathbf{A}}}$ , as both  $\hat{\mathbf{A}}$  and  $\mathbf{B}$  are full row rank. Then, the sufficient condition  $\text{Im}(\mathbf{B}) \subseteq \text{Im}(\hat{\mathbf{A}})$  holds.

The augmented Lagrangian function of (33) is formulated as

$$\mathcal{L}_{\rho, \epsilon}(\hat{\mathbf{x}}, \mathbf{y}, \boldsymbol{\lambda}) = \hat{f}(\hat{\mathbf{x}}) + h(\mathbf{y}) + \boldsymbol{\lambda}^\top (\hat{\mathbf{A}}\hat{\mathbf{x}} - \mathbf{B}\mathbf{y}) + \frac{\rho}{2} \|\hat{\mathbf{A}}\hat{\mathbf{x}} - \mathbf{B}\mathbf{y}\|_2^2 \quad (34)$$

The updates of the ADMM algorithm to optimize (33) are as follows

$$\begin{cases} \mathbf{y}^{k+1} &= \underset{\mathbf{y}}{\text{argmin}} \mathcal{L}_{\rho, \epsilon}(\hat{\mathbf{x}}^k, \mathbf{y}, \boldsymbol{\lambda}^k), \\ \hat{\mathbf{x}}^{k+1} &= \underset{\hat{\mathbf{x}}}{\text{argmin}} \mathcal{L}_{\rho, \epsilon}(\hat{\mathbf{x}}, \mathbf{y}^{k+1}, \boldsymbol{\lambda}^k), \\ \boldsymbol{\lambda}^{k+1} &= \boldsymbol{\lambda}^k + \rho(\hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}). \end{cases} \quad (35)$$

The optimality conditions of the variable sequence  $(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1})$  generated above are

$$\mathbf{B}^\top \boldsymbol{\lambda}^k + \rho \mathbf{B}^\top (\hat{\mathbf{A}}\hat{\mathbf{x}}^k - \mathbf{B}\mathbf{y}^{k+1}) = \mathbf{B}^\top \boldsymbol{\lambda}^{k+1} - \rho \mathbf{B}^\top \hat{\mathbf{A}}(\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k) \in \partial h(\mathbf{y}^{k+1}), \quad (36)$$

$$\nabla \hat{f}(\hat{\mathbf{x}}^{k+1}) + \hat{\mathbf{A}}^\top \boldsymbol{\lambda}^k + \rho \hat{\mathbf{A}}^\top (\hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}) = \nabla \hat{f}(\hat{\mathbf{x}}^{k+1}) + \hat{\mathbf{A}}^\top \boldsymbol{\lambda}^{k+1} = \mathbf{0}, \quad (37)$$

$$\frac{1}{\rho} (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k) = \hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}. \quad (38)$$

The convergence of this perturbed ADMM algorithm for the LS-LP problem is summarized in Theorem 2. The detailed proof is presented in the following sub-sections sequentially. Note that hereafter  $\|\cdot\|$  indicates the  $\ell_2$  norm for a vector, or the Frobenius norm for a matrix;  $\mathcal{A}_1 \succeq \mathcal{A}_2$  represents that  $\mathcal{A}_1 - \mathcal{A}_2$  is positive semi-definite, with  $\mathcal{A}_1, \mathcal{A}_2$  being square matrices;  $\nabla$  denotes the gradient operator,  $\nabla^2$  means the Hessian operator, and  $\partial$  is the sub-gradient operator;  $\mathbf{I}$  represents the identity matrix with compatible shape.

### A.1 Properties

In this section, we present some important properties of the objective function and constraints in (33), which will be used in the followed convergence analysis.

#### Properties on objective functions (P1)

- (P1.1)  $f, h$  and  $\mathcal{L}_{\rho, \epsilon}$  are semi-algebraic, lower semi-continuous functions and satisfy Kurdyka-Lojasiewicz (KL) property, and  $h$  is closed and proper
- (P1.2) There exist  $\mathcal{Q}_1, \mathcal{Q}_2$  such that  $\mathcal{Q}_1 \succeq \nabla^2 \hat{f}(\hat{\mathbf{x}}) \succeq \mathcal{Q}_2, \forall \hat{\mathbf{x}}$
- (P1.3)  $\liminf_{\|\hat{\mathbf{x}}\| \rightarrow \infty} \|\nabla \hat{f}(\hat{\mathbf{x}})\| = \infty$

#### Properties on constraint matrices (P2)

- (P2.1) There exists  $\sigma > 0$  such that  $\hat{\mathbf{A}}\hat{\mathbf{A}}^\top \succeq \sigma \mathbf{I}$

- (P2.2)  $\mathcal{Q}_2 + \rho \hat{\mathbf{A}}^\top \hat{\mathbf{A}} \succeq \delta \mathbf{I}$  for some  $\rho, \delta > 0$ , and  $\rho > \frac{1}{\epsilon}$
- (P2.3) There exists  $\mathcal{Q}_3 \succeq [\nabla^2 \hat{f}(\hat{\mathbf{x}})]^2, \forall \hat{\mathbf{x}}$ , and  $\delta \mathbf{I} \succ \frac{2}{\sigma \rho} \mathcal{Q}_3$
- (P2.4) Both  $\hat{\mathbf{A}}$  and  $\mathbf{B}$  are full row rank, and  $\text{Im}(\hat{\mathbf{A}}) \equiv \text{Im}(\mathbf{B}) \subseteq \mathbb{R}^{\text{rank of } \hat{\mathbf{A}}}$

**Remark. (1)** Although the definition of KL property (see Definition 1) is somewhat complex, but it holds for many widely used functions, according to [42]. Typical functions satisfying KL property includes: *a)* real analytic functions, and any polynomial function such as  $\|\mathbf{H}\mathbf{x} - \mathbf{b}\|$  belongs to this type; *b)* locally strongly convex functions, such as the logistic loss function  $\log(1 + \exp(-\mathbf{x}))$ ; *c)* semi-algebraic functions, such as  $\|\mathbf{x}\|_1, \|\mathbf{x}\|_2, \|\mathbf{x}\|_\infty, \|\mathbf{H}\mathbf{x} - \mathbf{b}\|_1, \|\mathbf{H}\mathbf{x} - \mathbf{b}\|_2, \|\mathbf{H}\mathbf{x} - \mathbf{b}\|_\infty$  and the indicator function  $\mathbb{I}(\cdot)$ . It is easy to verify that P1.1 holds in our problem. **(2)** Here we provide an instantiation of above hyper-parameters satisfying above properties. Firstly, it is easy to obtain that  $\nabla^2 \hat{f}(\hat{\mathbf{x}}) = \epsilon \mathbf{I}$ , and  $\hat{\mathbf{A}} \hat{\mathbf{A}}^\top = [\mathbf{A}, \epsilon \mathbf{I}][\mathbf{A}, \epsilon \mathbf{I}]^\top = \mathbf{A} \mathbf{A}^\top + \epsilon^2 \mathbf{I} \succ \epsilon^2 \mathbf{I}$ , as well as  $\rho \hat{\mathbf{A}}^\top \hat{\mathbf{A}} \succeq \epsilon \mathbf{I}$ , when  $\epsilon$  is small enough and  $\rho > \frac{1}{\epsilon}$  (e.g.,  $\rho = \frac{2}{\epsilon}$ ). Then, the values  $\mathcal{Q}_1 = \mathcal{Q}_2 = \epsilon \mathbf{I}, \mathcal{Q}_3 = \epsilon^2 \mathbf{I}, \delta = 2\epsilon, \sigma = \epsilon^2$  satisfy P1.2, P2.1, P2.2 and P2.3. Without loss of generality, we will adopt these specific values for these hyper-parameters to simplify the following analysis, while only keeping  $\rho$  and  $\epsilon$ .

## A.2 Decreasing of $\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)$

In this section, we firstly prove the decreasing property of the augmented Lagrangian function, *i.e.*,

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) > \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}), \forall k. \quad (39)$$

Firstly, utilizing P2.1, P2.3 and (37), we obtain that

$$\epsilon^2 \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|_2^2 \leq \|\hat{\mathbf{A}}(\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)\|_2^2 = \|\nabla \hat{f}(\hat{\mathbf{x}}^{k+1}) - \nabla \hat{f}(\hat{\mathbf{x}}^k)\|_2^2 = \epsilon^2 \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2. \quad (40)$$

Then, we have

$$\begin{aligned} & \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^k) = (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)^\top (\hat{\mathbf{A}} \hat{\mathbf{x}}^{k+1} - \mathbf{B} \mathbf{y}^{k+1}) \\ & = \frac{1}{\rho} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|_2^2 \leq \frac{1}{\rho} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 \end{aligned} \quad (41)$$

According to P1.2 and P2.2,  $\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}, \boldsymbol{\lambda}^k)$  is strongly convex with respect to  $\hat{\mathbf{x}}$ , with the parameter of at least  $2\epsilon$ . Then, we have

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^k) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \leq -\epsilon \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2. \quad (42)$$

As  $\mathbf{y}^{k+1}$  is the minimal solution of  $\mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)$ , it is easy to know

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \leq 0. \quad (43)$$

Combining (41), (42) and (43), we have

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \leq \left(\frac{1}{\rho} - \epsilon\right) \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 < 0, \quad (44)$$

where the last inequality utilizes P2.3 and  $\rho > \frac{1}{\epsilon}$ .

## A.3 Boundedness of $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$

Next, we prove the boundedness of  $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$ . We suppose that  $\rho$  is large enough such that there is  $0 < \gamma < \rho$  with

$$\inf_{\hat{\mathbf{x}}} \left( \hat{f}(\hat{\mathbf{x}}) - \frac{1}{2\epsilon^2 \gamma} \|\nabla \hat{f}(\hat{\mathbf{x}})\|_2^2 \right) = f^* > -\infty. \quad (45)$$

According to (44), for any  $k \geq 1$ , we have

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) = \hat{f}(\hat{\mathbf{x}}^k) + h(\mathbf{y}^k) + \frac{\rho}{2} \|\hat{\mathbf{A}} \hat{\mathbf{x}}^k - \mathbf{B} \mathbf{y}^k + \frac{\boldsymbol{\lambda}^k}{\rho}\|_2^2 - \frac{1}{2\rho} \|\boldsymbol{\lambda}^k\|_2^2 \leq \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^1, \hat{\mathbf{x}}^1, \boldsymbol{\lambda}^1) < \infty. \quad (46)$$

Besides, according to P2.1, we have

$$\epsilon^2 \|\boldsymbol{\lambda}^k\|_2^2 \leq \|\hat{\mathbf{A}}^\top \boldsymbol{\lambda}^k\|_2^2 = \|\nabla \hat{f}(\hat{\mathbf{x}}^k)\|_2^2. \quad (47)$$

Plugging (47) into (46), we obtain that

$$\infty > \hat{f}(\hat{\mathbf{x}}^k) + h(\mathbf{y}^k) + \frac{\rho}{2} \|\hat{\mathbf{A}} \hat{\mathbf{x}}^k - \mathbf{B} \mathbf{y}^k + \frac{\boldsymbol{\lambda}^k}{\rho}\|_2^2 - \frac{1}{2\epsilon^2 \rho} \|\nabla \hat{f}(\hat{\mathbf{x}}^k)\|_2^2 \geq f^* + \frac{\frac{1}{\rho} - \frac{1}{2\epsilon^2 \rho}}{2\epsilon^2} \|\nabla \hat{f}(\hat{\mathbf{x}}^k)\|_2^2 + h(\mathbf{y}^k) + \frac{\rho}{2} \|\hat{\mathbf{A}} \hat{\mathbf{x}}^k - \mathbf{B} \mathbf{y}^k + \frac{\boldsymbol{\lambda}^k}{\rho}\|_2^2. \quad (48)$$

According to the coerciveness of  $\nabla \hat{f}(\hat{\mathbf{x}}^k)$  (*i.e.*, P1.3), we obtain that  $\hat{\mathbf{x}}^k < \infty, \forall k$ , *i.e.*, the boundedness of  $\{\hat{\mathbf{x}}^k\}$ . From (47), we know the boundedness of  $\{\boldsymbol{\lambda}^k\}$ . Besides, according to P2.4,  $\{\hat{\mathbf{A}} \hat{\mathbf{x}}^k\}$  is also bounded. From (38), we obtain the boundedness of  $\{\mathbf{B} \mathbf{y}^k\}$ . Considering the full row rank of  $\mathbf{B}$  (*i.e.*, P2.4), the boundedness of  $\{\mathbf{y}^k\}$  is proved.

#### A.4 Convergence of Residual

According to the boundedness of  $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$ , there is a sub-sequence  $\{\mathbf{y}^{k_i}, \hat{\mathbf{x}}^{k_i}, \boldsymbol{\lambda}^{k_i}\}$  that converges to a cluster point  $\{\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*\}$ . Considering the lower semi-continuity of  $\mathcal{L}_{\rho, \epsilon}$  (i.e., P1.1), we have

$$\liminf_{i \rightarrow \infty} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k_i}, \hat{\mathbf{x}}^{k_i}, \boldsymbol{\lambda}^{k_i}) \geq \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*) > -\infty. \quad (49)$$

Summing (44) from  $k = M, \dots, N-1$  with  $M \geq 1$ , we have

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^M, \hat{\mathbf{x}}^M, \boldsymbol{\lambda}^M) \leq \left(\frac{1}{\rho} - \epsilon\right) \sum_{k=M}^{N-1} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 < 0. \quad (50)$$

Then, by setting  $N = k_i$  and  $M = 1$ , we have

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k_i}, \hat{\mathbf{x}}^{k_i}, \boldsymbol{\lambda}^{k_i}) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^1, \hat{\mathbf{x}}^1, \boldsymbol{\lambda}^1) \leq \left(\frac{1}{\rho} - \epsilon\right) \sum_{k=1}^{k_i-1} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2. \quad (51)$$

Taking limit on both sides of the above inequality, we obtain

$$-\infty < \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^1, \hat{\mathbf{x}}^1, \boldsymbol{\lambda}^1) \leq \left(\frac{1}{\rho} - \epsilon\right) \sum_{k=1}^{\infty} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 < 0. \quad (52)$$

It implies that

$$\lim_{k \rightarrow \infty} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| = 0. \quad (53)$$

Besides, according to (40), it is easy to obtain that

$$\lim_{k \rightarrow \infty} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\| = 0. \quad (54)$$

Moreover, utilizing  $\mathbf{B}\mathbf{y}^{k+1} = \hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \frac{1}{\rho}(\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)$  from (38), we have

$$\|\mathbf{B}(\mathbf{y}^{k+1} - \mathbf{y}^k)\| \leq \|\hat{\mathbf{A}}(\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k)\| + \frac{1}{\rho}\|(\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)\| + \frac{1}{\rho}\|(\boldsymbol{\lambda}^k - \boldsymbol{\lambda}^{k-1})\|. \quad (55)$$

Besides, as shown in Lemma 1 in [40], the full row rank of  $\mathbf{B}$  (i.e., P1.4) implies that

$$\|\mathbf{y}^{k+1} - \mathbf{y}^k\| \leq \bar{M}\|\mathbf{B}(\mathbf{y}^{k+1} - \mathbf{y}^k)\|, \quad (56)$$

where  $\bar{M} > 0$  is a constant. Taking limit on both sides of (55) and utilizing (56), we obtain

$$\lim_{k \rightarrow \infty} \|\mathbf{y}^{k+1} - \mathbf{y}^k\| = 0. \quad (57)$$

Combining (53), (54) and (57), we obtain that

$$\lim_{k \rightarrow \infty} \|\mathbf{y}^{k+1} - \mathbf{y}^k\|_2^2 + \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 + \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|_2^2 = 0. \quad (58)$$

By setting  $k+1 = k_i$ , plugging (53) into (36) and (54) into (37), and taking limit  $k_i \rightarrow \infty$ , we obtain the KKT conditions. It tells that the cluster point  $(\mathbf{y}^*, \hat{\mathbf{x}}^*)$  is the KKT point of LS-LP( $\theta; \epsilon$ ) (i.e., (33)).

#### A.5 Global Convergence

Inspired by the analysis presented in [25], in this section we will prove the following conclusions:

- $\sum_{k=1}^{\infty} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| < \infty$ ;
- $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$  converges to  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ ;
- $(\mathbf{y}^*, \hat{\mathbf{x}}^*)$  is the KKT point of (33).

Firstly, utilizing the optimality conditions (36, 37, 38), we have that

$$\partial_{\mathbf{y}} \mathcal{L}_{\rho, \epsilon}(\hat{\mathbf{x}}^{k+1}, \mathbf{y}^{k+1}, \boldsymbol{\lambda}^{k+1}) = \partial h(\mathbf{y}^{k+1}) - \mathbf{B}^T \boldsymbol{\lambda}^{k+1} - \rho \mathbf{B}^T (\hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}) \ni -\mathbf{B}^T (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k) - \rho \mathbf{B}^T \hat{\mathbf{A}}(\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k), \quad (59)$$

$$\nabla_{\hat{\mathbf{x}}} \mathcal{L}_{\rho, \epsilon}(\hat{\mathbf{x}}^{k+1}, \mathbf{y}^{k+1}, \boldsymbol{\lambda}^{k+1}) = \nabla_{\hat{\mathbf{x}}} \hat{f}(\hat{\mathbf{x}}^{k+1}) + \hat{\mathbf{A}}^T \boldsymbol{\lambda}^{k+1} + \rho \hat{\mathbf{A}}^T (\hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}) = \hat{\mathbf{A}}^T (\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k), \quad (60)$$

$$\nabla_{\boldsymbol{\lambda}} \mathcal{L}_{\rho, \epsilon}(\hat{\mathbf{x}}^{k+1}, \mathbf{y}^{k+1}, \boldsymbol{\lambda}^{k+1}) = \hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1} = \frac{1}{\rho}(\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k). \quad (61)$$

Further, combining with (40), there exists a constant  $C > 0$  such that

$$\text{dist}(0, \partial_{(\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda})} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1})) \leq C \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|, \quad (62)$$



where  $\text{dist}(\cdot, \cdot)$  denotes the distance between a vector and a set of vectors. Hereafter we denote  $\partial_{(\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda})} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1})$  as  $\partial \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1})$  for clarity. Besides, the relation (44) implies that there is a constant  $D \in (0, \epsilon - \frac{1}{\rho})$  such that

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}) \geq D \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2. \quad (63)$$

Moreover, the relation (49) implies that  $\{\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)\}$  is lower bounded along the convergent sub-sequence  $\{(\mathbf{y}^{k_i}, \hat{\mathbf{x}}^{k_i}, \boldsymbol{\lambda}^{k_i})\}$ . Combining with the its decreasing property, the limit of  $\{\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)\}$  exists. Thus, we will show that

$$\lim_{k \rightarrow \infty} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) = l^* := \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*). \quad (64)$$

To prove it, we utilize the fact that  $\mathbf{y}^{k+1}$  is the minimizer of  $\mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)$ , such that

$$\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \leq \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k). \quad (65)$$

Combining the above relation, (58) and the continuity of  $\mathcal{L}_{\rho, \epsilon}$  w.r.t.  $\hat{\mathbf{x}}$  and  $\boldsymbol{\lambda}$ , the following relation holds along the sub-sequence  $\{(\mathbf{y}^{k_i}, \hat{\mathbf{x}}^{k_i}, \boldsymbol{\lambda}^{k_i})\}$  that converges to  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ ,

$$\limsup_{i \rightarrow \infty} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k_i+1}, \hat{\mathbf{x}}^{k_i+1}, \boldsymbol{\lambda}^{k_i+1}) \leq \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*). \quad (66)$$

According to (58), the sub-sequence  $\{(\mathbf{y}^{k_i+1}, \hat{\mathbf{x}}^{k_i+1}, \boldsymbol{\lambda}^{k_i+1})\}$  also converges to  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ . Then, utilizing the lower semi-continuity of  $\mathcal{L}_{\rho, \epsilon}$ , we have

$$\liminf_{i \rightarrow \infty} \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k_i+1}, \hat{\mathbf{x}}^{k_i+1}, \boldsymbol{\lambda}^{k_i+1}) \geq \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*). \quad (67)$$

Combining (66) with (67), we know the existence of the limit of the sequence  $\{\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)\}$ , which proves the relation (64).

As  $\mathcal{L}_{\rho, \epsilon}$  is KL function, according to Definition 1, it has the following properties:

- There exist a constant  $\eta \in (0, \infty)$ , a continuous concave function  $\varphi : [0, \eta) \rightarrow \mathbb{R}_+$ , as well as a neighbourhood  $\mathcal{V}$  of  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ .  $\varphi$  is differentiable on  $(0, \eta)$  with positive derivatives.
- For all  $(\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda}) \in \mathcal{V}$  satisfying  $l^* < \mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda}) < l^* + \eta$ , we have

$$\varphi'(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda}) - l^*) \text{dist}(0, \partial \mathcal{L}_{\rho, \epsilon}(\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda})) \geq 1. \quad (68)$$

Then, we define the following neighborhood sets:

$$\mathcal{V}_\zeta := \left\{ (\mathbf{y}, \hat{\mathbf{x}}, \boldsymbol{\lambda}) \mid \|\hat{\mathbf{x}} - \hat{\mathbf{x}}^*\| < \zeta, \|\mathbf{y} - \mathbf{y}^*\| < \bar{M}(\|\hat{\mathbf{A}}\| + 1)\zeta, \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\| < \zeta \right\} \subseteq \mathcal{V} \quad (69)$$

$$\mathcal{V}_{\zeta, \hat{\mathbf{x}}} := \{ \hat{\mathbf{x}} \mid \|\hat{\mathbf{x}} - \hat{\mathbf{x}}^*\| < \zeta \}, \quad (70)$$

where  $\zeta > 0$  is a small constant.

Utilizing the relations (37) and (38), as well as P2.1, we obtain that for any  $k \geq 1$ , the following relation holds:

$$\epsilon^2 \|\boldsymbol{\lambda}^k - \boldsymbol{\lambda}^*\|_2^2 \leq \|\hat{\mathbf{A}}^\top (\boldsymbol{\lambda}^k - \boldsymbol{\lambda}^*)\|_2^2 = \|\nabla \hat{f}(\hat{\mathbf{x}}^k) - \nabla \hat{f}(\hat{\mathbf{x}}^*)\|_2^2 = \epsilon^2 \|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\|_2^2. \quad (71)$$

Also, the relations (37) and (38) imply that for any  $k \geq 1$ , we have

$$\|\mathbf{B}(\mathbf{y}^k - \mathbf{y}^*)\| = \|\hat{\mathbf{A}}(\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*) - \frac{1}{\rho}(\boldsymbol{\lambda}^k - \boldsymbol{\lambda}^{k-1})\| \leq \|\hat{\mathbf{A}}\| \|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| + \frac{1}{\rho} \|\boldsymbol{\lambda}^k - \boldsymbol{\lambda}^{k-1}\|. \quad (72)$$

Moreover, the relation (58) implies that  $\exists N_0 \geq 1$  such that  $\forall k \geq N_0$ , we have

$$\|\boldsymbol{\lambda}^k - \boldsymbol{\lambda}^{k-1}\| \leq \rho \zeta. \quad (73)$$

Similar to (56), the full row rank of  $\mathbf{B}$  implies  $\|\mathbf{y}^k - \mathbf{y}^*\| \leq \bar{M} \|\mathbf{B}(\mathbf{y}^k - \mathbf{y}^*)\|$ . Then, plugging (73) into (72), we obtain that

$$\|\mathbf{y}^k - \mathbf{y}^*\| \leq \bar{M}(\|\hat{\mathbf{A}}\| + 1)\zeta, \quad (74)$$

for any  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  and  $k \geq N_0$ . Combining (71) and (74), we know that if  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  and  $k \geq N_0$ , then  $(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \in \mathcal{V}_\zeta \subseteq \mathcal{V}$ .

Moreover, (44) and (64) implies that  $\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \geq l^*, \forall k \geq 1$ . Besides, as  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$  is a cluster point, we will obtain that  $\exists N \geq N_0$ , the following relations hold:

$$\begin{cases} \hat{\mathbf{x}}^N \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}} \\ l^* < \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) < l^* + \eta \\ \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + 2\sqrt{(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) - l^*)/D} + \frac{C}{D}(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) - l^*) < \zeta \end{cases} \quad (75)$$

Next, We will show that if  $\hat{\mathbf{x}}^N \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  and  $l^* < \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) < l^* + \eta$  hold for some fixed  $k \geq N_0$ , then the following relation holds

$$\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| + (\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| - \|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^{k-1}\|) \leq \frac{C}{D} \left[ \varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - l^*) - \varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}) - l^*) \right]. \quad (76)$$

To prove (76), we utilize the fact that  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$ ,  $k \geq N_0$  implies that  $(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}} \subseteq \mathcal{V}$ . And, combining with  $l^* < \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) < l^* + \eta$ , we obtain that

$$\varphi'(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - l^*) \text{dist}(0, \partial \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)) \geq 1. \quad (77)$$

Combining the relations (62), (63) and (77), as well as the concavity of  $\varphi$ , we obtain that

$$\begin{aligned} & C \|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^{k-1}\| \cdot [\varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - l^*) - \varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}) - l^*)] \\ & \geq \text{dist}(0, \partial \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)) \cdot [\varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - l^*) - \varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1}) - l^*)] \\ & \geq \text{dist}(0, \partial \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)) \cdot \varphi'(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - l^*) \cdot [\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{k+1}, \hat{\mathbf{x}}^{k+1}, \boldsymbol{\lambda}^{k+1})] \\ & \geq D \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|^2, \end{aligned} \quad (78)$$

for all such  $k$ . Taking square root on both sides of (78), and utilizing the fact that  $a + b \geq 2\sqrt{ab}$ , then (76) is proved.

We then prove  $\forall k \geq N$ ,  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  holds. This claim can be proved through induction. Obviously it is true for  $k = N$  by construction, as shown in (75). For  $k = N + 1$ , we have

$$\begin{aligned} & \|\hat{\mathbf{x}}^{N+1} - \hat{\mathbf{x}}^*\| \leq \|\hat{\mathbf{x}}^{N+1} - \hat{\mathbf{x}}^N\| + \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| \\ & \leq \sqrt{(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) - \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{N+1}, \hat{\mathbf{x}}^{N+1}, \boldsymbol{\lambda}^{N+1}))/D} + \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| \\ & \leq \sqrt{(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N) - l^*)/D} + \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| < \zeta, \end{aligned} \quad (79)$$

where the first inequality utilizes (63), and the last inequality follows the last relation in (75). Thus,  $\forall k \geq N$ ,  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  holds.

Next, we suppose that  $\hat{\mathbf{x}}^N, \dots, \hat{\mathbf{x}}^{N+t-1} \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  for some  $t > 1$ , and we need to prove that  $\hat{\mathbf{x}}^{N+t} \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  also holds, *i.e.*,

$$\begin{aligned} & \|\hat{\mathbf{x}}^{N+t} - \hat{\mathbf{x}}^*\| \leq \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + \|\hat{\mathbf{x}}^{N+1} - \hat{\mathbf{x}}^N\| + \sum_{i=1}^{t-1} \|\hat{\mathbf{x}}^{N+i+1} - \hat{\mathbf{x}}^{N+i}\| \\ & = \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + 2\|\hat{\mathbf{x}}^{N+1} - \hat{\mathbf{x}}^N\| - \|\hat{\mathbf{x}}^{N+t} - \hat{\mathbf{x}}^{N+t-1}\| + \sum_{i=1}^{t-1} \left[ \|\hat{\mathbf{x}}^{N+i+1} - \hat{\mathbf{x}}^{N+i}\| + (\|\hat{\mathbf{x}}^{N+i+1} - \hat{\mathbf{x}}^{N+i}\| - \|\hat{\mathbf{x}}^{N+i} - \hat{\mathbf{x}}^{N+i-1}\|) \right] \\ & \leq \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + 2\|\hat{\mathbf{x}}^{N+1} - \hat{\mathbf{x}}^N\| + \frac{C}{D} \sum_{i=1}^{t-1} [\varphi^{N+i} - \varphi^{N+i+1}] \\ & \leq \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + 2\|\hat{\mathbf{x}}^{N+1} - \hat{\mathbf{x}}^N\| + \frac{C}{D} \sum_{i=1}^{t-1} \varphi^{N+1} \\ & \leq \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + 2\sqrt{\frac{\mathcal{L}_{\rho, \epsilon}^N - \mathcal{L}_{\rho, \epsilon}^{N+1}}{D}} + \frac{C}{D} \sum_{i=1}^{t-1} \varphi^{N+1} \\ & \leq \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^*\| + 2\sqrt{\frac{\mathcal{L}_{\rho, \epsilon}^N - l^*}{D}} + \frac{C}{D} \sum_{i=1}^{t-1} \varphi^{N+1} < \zeta \end{aligned} \quad (80)$$

where  $\varphi^{N+i} = \varphi(\mathcal{L}_{\rho, \epsilon}(\mathbf{y}^{N+i}, \hat{\mathbf{x}}^{N+i}, \boldsymbol{\lambda}^{N+i}) - l^*)$  and  $\mathcal{L}_{\rho, \epsilon}^N = \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^N, \hat{\mathbf{x}}^N, \boldsymbol{\lambda}^N)$ . The second inequality follows from (76). The fourth inequality follows from (63). The fifth inequality utilizes the fact that  $\mathcal{L}_{\rho, \epsilon}^{N+1} > l^*$ , and the last inequality follows from the last relation in (75). Thus,  $\hat{\mathbf{x}}^{N+k} \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  holds. We have proved that  $\forall k \geq N$ ,  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$  holds by induction.

Then, according to  $\forall k \geq N$ ,  $\hat{\mathbf{x}}^k \in \mathcal{V}_{\zeta, \hat{\mathbf{x}}}$ , we can sum both sides of (76) from  $k = N$  to  $\infty$ , to obtain that

$$\sum_{k=N}^{\infty} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| \leq \frac{C}{D} \varphi^N + \|\hat{\mathbf{x}}^N - \hat{\mathbf{x}}^{N-1}\| < \infty, \quad (81)$$

which implies that  $\sum_{k=1}^{\infty} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| < \infty$  holds. Thus  $\{\hat{\mathbf{x}}^k\}$  converges. The convergence of  $\{\mathbf{y}^k\}$  follows from  $\mathbf{B}\mathbf{y}^{k+1} = \hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} + \frac{1}{\rho}(\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k)$  in (38) and (58), as well as the surjectivity of  $\mathbf{B}$  (*i.e.*, full row rank). The convergence of  $\{\boldsymbol{\lambda}^k\}$  follows from  $\nabla \hat{f}(\hat{\mathbf{x}}^{k+1}) = -\hat{\mathbf{A}}^\top \boldsymbol{\lambda}^{k+1}$  in (37) and the surjectivity of  $\hat{\mathbf{A}}$  (*i.e.*, full row rank). Consequently,  $\{\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k\}$  converges to the cluster point  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$ . The conclusion that  $(\mathbf{y}^*, \hat{\mathbf{x}}^*)$  is the KKT point of Problem (33) has been proved in Section A.4.

## A.6 $\epsilon$ -KKT Point of the Original LS-LP Problem

**Proposition 1** *The globally converged solution  $(\mathbf{y}^*, \mathbf{x}^*, \boldsymbol{\lambda}^*)$  produced by the ADMM algorithm for the perturbed LS-LP problem (33) is the  $\epsilon$ -KKT solution to the original LS-LP problem (32).*

*Proof* The globally converged solution  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$  to the perturbed LS-LP problem (33) satisfies the following relations:

$$\mathbf{B}^\top \boldsymbol{\lambda}^* \in \partial h(\mathbf{y}^*), \quad \nabla \hat{f}(\hat{\mathbf{x}}^*) = -\hat{\mathbf{A}}^\top \boldsymbol{\lambda}^*, \quad \hat{\mathbf{A}}\hat{\mathbf{x}}^* = \mathbf{B}\mathbf{y}^*. \quad (82)$$

Recalling the definitions  $\hat{\mathbf{A}} = [\mathbf{A}, \epsilon \mathbf{I}]$ ,  $\hat{\mathbf{x}} = [\mathbf{x}; \bar{\mathbf{x}}]$  and  $\hat{f}(\hat{\mathbf{x}}) = f(\mathbf{x}) + \frac{\epsilon}{2} \hat{\mathbf{x}}^\top \hat{\mathbf{x}}$ , the above relations imply that

$$\nabla \hat{f}(\hat{\mathbf{x}}^*) + \hat{\mathbf{A}}^\top \boldsymbol{\lambda}^* = \nabla f(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\lambda}^* + \epsilon \mathbf{x}^* = \mathbf{0} \Rightarrow \|\nabla f(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\lambda}^*\| = \epsilon \|\mathbf{x}^*\| = O(\epsilon), \quad (83)$$

$$\hat{\mathbf{A}} \hat{\mathbf{x}}^* + \mathbf{B} \mathbf{y}^* = \mathbf{A} \mathbf{x}^* + \epsilon \bar{\mathbf{x}}^* + \mathbf{B} \mathbf{y}^* = \mathbf{0} \Rightarrow \|\mathbf{A} \mathbf{x}^* + \mathbf{B} \mathbf{y}^*\| = \|\epsilon \bar{\mathbf{x}}^*\| = O(\epsilon), \quad (84)$$

where we utilize the boundedness of  $\{\hat{\mathbf{x}}^*\}$ . Thus, according to Definition 2, the globally converged point  $(\mathbf{y}^*, \mathbf{x}^*)$  is the  $\epsilon$ -KKT solution to the original LS-LP problem (32).

## A.7 Convergence Rate

**Lemma 3** *Firstly, without loss of generality, we can assume that  $l^* = \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*) = 0$  (e.g., one can replace  $l_k = \mathcal{L}_{\rho, \epsilon}(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)$  by  $l_k - l^*$ ). We further assume that  $\mathcal{L}_{\rho, \epsilon}$  has the KL property at  $(\mathbf{y}^*, \hat{\mathbf{x}}^*, \boldsymbol{\lambda}^*)$  with the concave function  $\varphi(s) = cs^{1-p}$ , where  $p \in [0, 1)$ ,  $c > 0$ . Consequently, we can obtain the following inequalities:*

- (i) if  $p = 0$ , then  $\{(\mathbf{y}^k, \hat{\mathbf{x}}^k, \boldsymbol{\lambda}^k)\}_{k=1, \dots, \infty}$  can converge in finite steps;
- (ii) If  $p \in (0, \frac{1}{2}]$ , then there exist  $c > 0$  and  $\tau \in (0, 1)$  such that  $\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| \leq c\tau^k$ ;
- (iii)  $p \in (\frac{1}{2}, 1)$ , then there exist  $c > 0$  such that  $\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| \leq ck^{-\frac{1-p}{2p-1}}$ .

*Proof* (i) If  $p = 0$ , we define a subset  $H = \{k \in \mathbb{N} : \hat{\mathbf{x}}_k \neq \hat{\mathbf{x}}_{k+1}\}$ . If  $k \in H$  is sufficiently large, then there exists  $C_3 > 0$  such that

$$\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|^2 \geq C_3 > 0. \quad (85)$$

Combining with (63), we have

$$l_k - l_{k+1} \geq D \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|^2 \geq C_3 D > 0. \quad (86)$$

If the subset  $H$  is infinite, then it will contradict to the fact that  $l_k - l_{k+1} \rightarrow 0$  as  $k \rightarrow \infty$ . Thus,  $H$  is a finite subset, leading to the conclusion that  $\{\hat{\mathbf{x}}^k\}_{k \in \mathbb{N}}$  will converge in finite steps. Recalling the relationships between  $\hat{\mathbf{x}}_k$  and  $\mathbf{y}_k, \boldsymbol{\lambda}_k$  (see the descriptions under (81)), we also obtain that  $\{\mathbf{y}_k, \boldsymbol{\lambda}_k\}_{k \in \mathbb{N}}$  converges in finite steps.

By defining  $\Delta_k = \sum_k^\infty \|x^k + 1 - x^k\|$ , the inequality (81) can be rewritten as follows

$$\Delta_k \leq \frac{C}{D} \varphi(l_k) + (\Delta_{k-1} - \Delta_k) < \infty. \quad (87)$$

Besides, the KL property and  $l^* = 0$  give that

$$\varphi'(l_k) \text{dist}(0, \partial(l_k)) = c(1-p)l_k^{1-p} \text{dist}(0, \partial(l_k)) \geq 1 \Rightarrow l_k^p \leq c(1-p) \text{dist}(0, \partial(l_k)). \quad (88)$$

Combining with (62), we obtain

$$l_k^p \leq c(1-p)C(\Delta_{k-1} - \Delta_k) \Rightarrow \varphi(l_k) = cl_k^{1-p} \leq c(c(1-p)C)^{\frac{1-p}{p}} (\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}} = C_1 (\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}}. \quad (89)$$

Then, inserting (89) into (87), we obtain

$$\Delta_k \leq C_2 (\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}} + (\Delta_{k-1} - \Delta_k) < \infty. \quad (90)$$

(ii) If  $p \in (0, \frac{1}{2}]$ , then  $\frac{1-p}{p} \geq 1$ . Besides, since  $(\Delta_{k-1} - \Delta_k) \rightarrow 0$  when  $k \rightarrow \infty$ , there exists an integer  $K_0$  such that  $(\Delta_{k-1} - \Delta_k) < 1$ , leading to that  $(\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}} \leq (\Delta_{k-1} - \Delta_k)$ . Inserting it into (90), we obtain that

$$\Delta_k \leq (C_2 + 1)(\Delta_{k-1} - \Delta_k) \Rightarrow \Delta_k \leq C_3 (\Delta_{k-1} - \Delta_k) \Rightarrow \Delta_k \leq \frac{C_3}{1 + C_3} \Delta_{k-1} = \tau \Delta_{k-1}, \text{ with } \tau \in (0, 1), \forall k > K_0. \quad (91)$$

It is easy to deduce that  $\Delta_k \leq (\Delta_{K_0} \tau^{-K_0}) \tau^k = \frac{c}{2} \tau^k$ , with  $c$  being a positive constant. Note that  $k$  in  $\tau^k$  indicates  $k$  power of  $\tau$ , rather than the iteration index. Combining with  $\|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| \leq \Delta_k$ , it is easy to obtain that  $\|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| \leq \frac{c}{2} \tau^k$  with  $\tau \in (0, 1)$  and  $c$  being a positive constant. Then, we have

$$\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| \leq \|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| + \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^*\| \leq \frac{c}{2} (\tau^{k+1} + \tau^k) \leq c\tau^k. \quad (92)$$

(iii) If  $p \in (\frac{1}{2}, 1)$ , then  $\frac{1-p}{p} < 1$ . Then, it is easy to obtain that  $(\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}} > (\Delta_{k-1} - \Delta_k)$ . Inserting it into (90), we obtain that

$$\Delta_k \leq (C_2 + 1)(\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}} \Rightarrow \Delta_k \leq C_3 (\Delta_{k-1} - \Delta_k)^{\frac{1-p}{p}} \Rightarrow \Delta_k^{\frac{p}{1-p}} \leq C_4 (\Delta_{k-1} - \Delta_k), \forall k > K_0. \quad (93)$$

It has been studied in Theorem 2 of [1] that the above inequality can deduce  $\Delta_k \leq \frac{c}{2} k^{-\frac{1-p}{2p-1}}$ , with  $c$  being a positive constant. Since  $\|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| \leq \Delta_k$ , we have that  $\|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| \leq \frac{c}{2} k^{-\frac{1-p}{2p-1}}$ . Then, we have

$$\|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\| \leq \|\hat{\mathbf{x}}^k - \hat{\mathbf{x}}^*\| + \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^*\| \leq \frac{c}{2} (k^{-\frac{1-p}{2p-1}} + (k+1)^{-\frac{1-p}{2p-1}}) \leq ck^{-\frac{1-p}{2p-1}}. \quad (94)$$

**Proposition 2** We adopt the same assumptions in Lemma 3. Then,

- (i) If  $p = 0$ , then we will obtain the  $\epsilon$ -KKT solution to the LS-LP problem in finite steps.
- (ii) If  $p \in (0, \frac{1}{2}]$ , then we will obtain the  $\epsilon$ -KKT solution to the LS-LP problem in at least  $O(\log_{\frac{1}{\tau}}(\frac{1}{\epsilon})^2)$  steps.
- (iii) If  $p \in (\frac{1}{2}, 1)$ , then we will obtain the  $\epsilon$ -KKT solution to the LS-LP problem in at least  $O((\frac{1}{\epsilon})^{\frac{4p-2}{1-p}})$  steps.

*Proof* The conclusion (i) directly holds from Lemma 3(i).

According to the optimality condition (36), we have

$$\begin{aligned} \text{dist}(\mathbf{B}^\top \boldsymbol{\lambda}^{k+1}, \partial h(\mathbf{y}^{k+1})) &= \|\rho \mathbf{B}^\top \hat{\mathbf{A}}(\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k)\|_2 \\ \Rightarrow \text{dist}^2(\mathbf{B}^\top \boldsymbol{\lambda}^{k+1}, \partial h(\mathbf{y}^{k+1})) &= \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_{\rho^2 \hat{\mathbf{A}}^\top \mathbf{B} \mathbf{B}^\top \hat{\mathbf{A}}}^2 \leq \xi_{\max}(\rho^2 \hat{\mathbf{A}}^\top \mathbf{B} \mathbf{B}^\top \hat{\mathbf{A}}) \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 = O(\frac{1}{\epsilon}) \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2^2 \\ \Rightarrow \text{dist}(\mathbf{B}^\top \boldsymbol{\lambda}^{k+1}, \partial h(\mathbf{y}^{k+1})) &\leq O(\frac{1}{\epsilon}) \cdot \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2 \end{aligned} \quad (95)$$

According to the optimality condition (38) and the relation (40), we obtain that

$$\|\hat{\mathbf{A}}\hat{\mathbf{x}}^{k+1} - \mathbf{B}\mathbf{y}^{k+1}\|_2 = \frac{1}{\rho} \|\boldsymbol{\lambda}^{k+1} - \boldsymbol{\lambda}^k\|_2 \leq \frac{1}{\rho} \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2 \leq \epsilon \cdot \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2 \leq O(\frac{1}{\epsilon}) \cdot \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2. \quad (96)$$

According to Lemma 3, we have

- (ii) If  $p \in (0, \frac{1}{2}]$ , then

$$O(\frac{1}{\epsilon}) \cdot \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2 \leq O(\frac{1}{\epsilon}) \tau^k \leq O(\epsilon) \Rightarrow k \geq O(\log_{\frac{1}{\tau}}(\frac{1}{\epsilon})^2),$$

which means that when  $k \geq O(\log_{\frac{1}{\tau}}(\frac{1}{\epsilon})^2)$ , we will obtain the  $\epsilon$ -KKT solution to the perturbed LS-LP problem, *i.e.*, the  $\epsilon$ -KKT solution to the original LS-LP problem.

- (iii) If  $p \in (\frac{1}{2}, 1)$ , then

$$O(\frac{1}{\epsilon}) \cdot \|\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k\|_2 \leq O(\frac{1}{\epsilon}) k^{-\frac{1-p}{2p-1}} \leq O(\epsilon) \Rightarrow k \geq O((\frac{1}{\epsilon})^{\frac{4p-2}{1-p}}),$$

which means that when  $k \geq O((\frac{1}{\epsilon})^{\frac{4p-2}{1-p}})$ , we will obtain the  $\epsilon$ -KKT solution to the perturbed LS-LP problem, *i.e.*, the  $\epsilon$ -KKT solution to the original LS-LP problem.

## References

1. Attouch, H., Bolte, J.: On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming* **116**(1-2), 5–16 (2009)
2. Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the kurdyka-lojasiewicz inequality. *Mathematics of Operations Research* **35**(2), 438–457 (2010)
3. Besag, J.: On the statistical analysis of dirty pictures. *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 259–302 (1986)
4. Bolte, J., Daniilidis, A., Lewis, A.: The lojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization* **17**(4), 1205–1223 (2007)
5. Bolte, J., Daniilidis, A., Lewis, A., Shiota, M.: Clarke subgradients of stratifiable functions. *SIAM Journal on Optimization* **18**(2), 556–572 (2007)
6. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning* **3**(1), 1–122 (2011)
7. Elidan, G., Globerson, A., Heinemann, U.: Pascal 2011 probabilistic inference challenge. <http://www.cs.huji.ac.il/project/PASCAL/index.php> (2012)
8. Fu, Q., Banerjee, H.W.A.: Bethe-admm for tree decomposition based parallel map inference. In: *Uncertainty in Artificial Intelligence*, p. 222. Citeseer (2013)
9. Globerson, A., Jaakkola, T.S.: Fixing max-product: Convergent message passing algorithms for map lp-relaxations. In: *NIPS*, pp. 553–560 (2008)
10. Gould, S., Fulton, R., Koller, D.: Decomposing a scene into geometric and semantically consistent regions. In: *ICCV*, pp. 1–8. IEEE (2009)
11. Jaimovich, A., Elidan, G., Margalit, H., Friedman, N.: Towards an integrated protein–protein interaction network: A relational markov network approach. *Journal of Computational Biology* **13**(2), 145–164 (2006)
12. Johnson, J.K., Malioutov, D.M., Willsky, A.S.: Lagrangian relaxation for map estimation in graphical models. *arXiv preprint arXiv:0710.0013* (2007)
13. Jojic, V., Gould, S., Koller, D.: Accelerated dual decomposition for map inference. In: *ICML*, pp. 503–510 (2010)
14. Kappes, J.H., Andres, B., Hamprecht, F.A., Schnörr, C., Nowozin, S., Batra, D., Kim, S., Kausler, B.X., Kröger, T., Lellmann, J., et al.: A comparative study of modern inference techniques for structured discrete energy minimization problems. *International Journal of Computer Vision* **115**, 155–184 (2015)
15. Kappes, J.H., Savchynskyy, B., Schnörr, C.: A bundle approach to efficient map-inference by lagrangian relaxation. In: *CVPR*, pp. 1688–1695. IEEE (2012)
16. Karush, W.: Minima of functions of several variables with inequalities as side constraints. M. Sc. Dissertation. Dept. of Mathematics, Univ. of Chicago (1939)
17. Kelley, J.: The cutting-plane method for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics* pp. 703–712 (1960)
18. Koller, D., Nir, F. (eds.): *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, Cambridge, MA (2009)
19. Kolmogorov, V.: Convergent tree-reweighted message passing for energy minimization. *IEEE transactions on pattern analysis and machine intelligence* **28**(10), 1568–1583 (2006)
20. Komodakis, N., Paragios, N., Tziritas, G.: Mrf optimization via dual decomposition: Message-passing revisited. In: *ICCV*, pp. 1–8. IEEE (2007)
21. Kschischang, F.R., Frey, B.J., Loeliger, H.A.: Factor graphs and the sum-product algorithm. *IEEE Transactions on information theory* **47**(2), 498–519 (2001)
22. Kuhn, H.W., Tucker, A.W.: Nonlinear programming. In: *Traces and emergence of nonlinear programming*, pp. 247–258. Springer (2014)
23. Land, A.H., Doig, A.G.: An automatic method of solving discrete programming problems. *Econometrica: Journal of the Econometric Society* pp. 497–520 (1960)
24. Laurent, M., Rendl, F.: *Semidefinite programming and integer programming*. Centrum voor Wiskunde en Informatica (2002)
25. Li, G., Pong, T.K.: Global convergence of splitting methods for nonconvex composite optimization. *SIAM Journal on Optimization* **25**(4), 2434–2460 (2015)
26. Lojasiewicz, S.: Une propriété topologique des sous-ensembles analytiques réels. *Les équations aux dérivées partielles* **117**, 87–89 (1963)
27. Martins, A.F., Figueiredo, M.A., Aguiar, P.M., Smith, N.A., Xing, E.P.: An augmented lagrangian approach to constrained map inference. In: *ICML* (2011)
28. Martins, A.F., Figueiredo, M.A., Aguiar, P.M., Smith, N.A., Xing, E.P.: Ad3: Alternating directions dual decomposition for map inference in graphical models. *Journal of Machine Learning Research* **16**, 495–545 (2015)
29. Meshi, O., Globerson, A.: An alternating direction method for dual map lp relaxation. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 470–483. Springer (2011)
30. Meshi, O., Mahdavi, M., Schwing, A.: Smooth and strong: Map inference with linear convergence. In: *NIPS*, pp. 298–306 (2015)
31. Otten, L., Dechter, R.: Anytime and/or depth-first search for combinatorial optimization. *AI Communications* **25**(3), 211–227 (2012)
32. Otten, L., Ihler, A., Kask, K., Dechter, R.: Winning the pascal 2011 map challenge with enhanced and/or branch-and-bound. In: *IN NIPS WORKSHOP DISCML*. Citeseer (2012)
33. Savchynskyy, B., Schmidt, S., Kappes, J., Schnörr, C.: Efficient mrf energy minimization via adaptive diminishing smoothing. *arXiv preprint arXiv:1210.4906* (2012)
34. Schwing, A.G., Hazan, T., Pollefeys, M., Urtasun, R.: Globally convergent dual map lp relaxation solvers using fenchel-young margins. In: *NIPS*, pp. 2384–2392 (2012)
35. Schwing, A.G., Hazan, T., Pollefeys, M., Urtasun, R.: Globally convergent parallel map lp relaxation solver using the frank-wolfe algorithm. In: *ICML*, pp. 487–495 (2014)
36. Sontag, D.A.: Approximate inference in graphical models using lp relaxations. Ph.D. thesis, Massachusetts Institute of Technology (2010)
37. Sontag, D.A., Li, Y., et al.: Efficiently searching for frustrated cycles in map inference. In: *UAI* (2012)
38. Wainwright, M.J., Jaakkola, T.S., Willsky, A.S.: Map estimation via agreement on trees: message-passing and linear programming. *IEEE transactions on information theory* **51**(11), 3697–3717 (2005)
39. Wainwright, M.J., Jordan, M.I.: Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning* **1**(1-2), 1–305 (2008)
40. Wang, Y., Yin, W., Zeng, J.: Global convergence of admm in non-convex nonsmooth optimization. *Journal of scientific programming* (2017)
41. Wu, B., Ghanem, B.:  $\ell_p$ -box admm: A versatile framework for integer programming. *IEEE transactions on pattern analysis and machine intelligence* **41**(7), 1695–1708 (2019)
42. Xu, Y., Yin, W.: A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences* **6**(3), 1758–1789 (2013)