

# Spurious Stationarity and Hardness Results for Mirror Descent

**He Chen**

*The Chinese University of Hong Kong*

HCHEN@SE.CUHK.EDU.HK

**Jiajin Li**

*The University of British Columbia*

JIAJIN.LI@SAUDER.UBC.CA

**Anthony Man-Cho So**

*The Chinese University of Hong Kong*

MANCHOSO@SE.CUHK.EDU.HK

*Authors ordered alphabetically*

## Abstract

Despite the success of Bregman proximal-type algorithms, such as mirror descent, in machine learning, most theoretical results depend on the gradient Lipschitz property of the kernel, excluding widely used cases like the Shannon entropy kernel. This paper uncovers a fundamental limitation: *Spurious stationary points* inevitably arise when non-gradient Lipschitz kernels are used. We establish an algorithm-dependent hardness result, showing that Bregman proximal-type algorithms cannot escape these spurious stationary points in finite steps if the initial point is unfavorable, even in convex settings. Those challenges are discovered through the lack of a well-defined stationarity measure, typically based on Bregman divergence, for these algorithms. While some extensions attempt to address this, we demonstrate that they still fail to distinguish reliably between stationary and non-stationary points for non-gradient Lipschitz kernels. Our findings highlight the need for new theoretical tools and algorithms within Bregman geometry, opening new avenues for further research.

## 1. Introduction

In this paper, we study structured nonsmooth (non)-convex optimization problems of the form

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}), \quad (\text{P})$$

where  $\text{dom}(g) = \mathcal{X}$  is a nonempty closed convex set,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a continuously differentiable function, and  $g : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$  is a convex and locally Lipschitz continuous function. To solve Problem (P) efficiently, Bregman proximal-type (BPs) algorithms are widely used methods that effectively exploit the geometry of the set  $\mathcal{X}$  to avoid costly projection or proximal steps under the Euclidean space [1, 6, 8, 26].

We now present the unified formulation of the class of Bregman proximal-type (BPs) algorithms that we will investigate in this paper, where the iterative schemes are defined as follows

$$\mathbf{x}^{k+1} = \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^n} \left\{ \gamma(\mathbf{y}; \mathbf{x}^k) + g(\mathbf{y}) + \frac{1}{t_k} D_h(\mathbf{y}, \mathbf{x}^k) \right\}. \quad (1.1)$$

Here,  $\gamma(\cdot; \mathbf{x})$  is the surrogate model for  $f$  at the point  $\mathbf{x}$ ,  $t_k \geq 0$  is the step size in the  $k$ -th iteration and  $D_h$  is the Bregman divergence associated with the kernel function  $h$ . Specifically, when the

surrogate model is the original function  $f$ , the update rule (1.1) simplifies to the Bregman proximal point method [11, 17]. If the surrogate model is a linear approximation, given by  $\gamma(\mathbf{y}; \mathbf{x}^k) := f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T(\mathbf{x} - \mathbf{x}^k)$ , then the update (1.1) encompasses the Bregman proximal gradient descent method (BPG) [3, 5, 10, 24, 29]. Alternatively, one may choose a quadratic surrogate model given by  $\gamma(\mathbf{y}; \mathbf{x}^k) := f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T(\mathbf{x} - \mathbf{x}^k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^k)^T \nabla^2 f(\mathbf{x}^k)(\mathbf{x} - \mathbf{x}^k)$ , which has been recently investigated by Doikov and Nesterov [13].

Many efforts have been devoted to obtaining non-asymptotic convergence results for Bregman proximal-type algorithms, which are expected to be similar to those for Euclidean cases. However, these results have been limited to scenarios where  $\text{dom}(h) = \mathbb{R}^n$  (see, e.g., [27]) or where  $\nabla h$  is Lipschitz continuous (see, e.g., [28]). These conditions essentially guarantee the non-degenerate property of the mirror map, where Euclidean and Bregman geometries align. Unfortunately, they exclude most of powerful kernels that enhance the practical appeal of BPs, such as the Shannon entropy function. This raises a natural question: *Is there any fundamental difference between Bregman proximal-type algorithms and their Euclidean counterparts when the mirror map is degenerate (i.e., when the kernel function is not gradient Lipschitz)?*

In this paper, we fully address this question and explore the fundamental differences that prevent achieving analogous convergence results for non-gradient Lipschitz kernels. We begin by identifying a class of undesirable points, termed *spurious stationary points*, which arise in scenarios where the kernel function is not gradient Lipschitz continuous. We then show that the existence of spurious stationary points is inevitable, even in simple convex problems. More importantly, if we initialize BPs near a spurious stationary point, the generated sequence can remain trapped within a small neighborhood of this point, regardless of the number of steps taken. This hardness result provides a negative answer to the possibility of obtaining similar convergence results as those in the Euclidean space.

Taking a step further, we aim to clarify why BPs become trapped at these spurious stationary points and how we discovered this result that has been overlooked by the research community for an extended period. To analyze the behavior of the iterate sequences and their proximity to stationary points, researchers typically introduce a residual function  $R : \mathbb{R}^n \rightarrow \mathbb{R}_+$  quantifies the stationarity of the iterates and subsequently establish the convergence of  $\{R(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ . This trapping phenomenon is observed through the unsatisfactory behaviours of all existing stationarity measures. Specifically, these measures fail to satisfy the following equivalence (Q) when the kernel function is not gradient Lipschitz continuous:

$$\boxed{\lim_{k \rightarrow \infty} R(\mathbf{x}^k) = 0 \stackrel{?}{\iff} \mathbf{0} \in \partial F \left( \lim_{k \rightarrow \infty} \mathbf{x}^k \right).} \quad (\text{Q})$$

This failure implies that we cannot classify whether the iterates of BPs are approximately stationary or not, even when the stationarity measures are relatively small. Spurious stationary points are precisely those points where the stationarity measure equals zero, yet zero does not belong to their subdifferential. It is worth noting that the equivalence (Q) trivially holds for all algorithms under Euclidean geometry, as well as under conditions such as  $\text{dom}(h) = \mathbb{R}^n$  (see, e.g., [27]) or when  $\nabla h$  is Lipschitz continuous (see, e.g., [28]) for BPs. Consequently, there are no spurious stationary points in these cases.

**Main technical novelty** The key tool for establishing the existence of spurious stationary points and the associated hardness results is the *extended Bregman stationarity measure*, which may also

be of independent interest. We begin by unifying all existing stationarity measures and extending their definitions to encompass the entire domain of kernels. This extension allows us to compute the residual function at accumulation points. Notably, prior to our work, all existing stationarity measures were not well-defined at the boundaries of the kernel functions' domains. We then establish the continuity properties of the newly developed stationarity measure and demonstrate that a stationarity measure equal to zero can still indicate non-stationary points (i.e., spurious stationary points). The hardness results follow directly from this analysis.

**Notation.** Our notation is mostly standard. The set of extended real numbers is denoted by  $\overline{\mathbb{R}}$ . For a vector  $\mathbf{x} \in \mathbb{R}^n$ , its  $i$ -th coordinate is represented by  $x_i$ , and  $\mathbf{x}_{\mathcal{I}}$  denotes a subvector of  $\mathbf{x}$  indexed by  $\mathcal{I}$ . The Euclidean ball  $\mathbb{B}_\epsilon(\mathbf{x})$  is defined through  $\mathbb{B}_\epsilon(\mathbf{x}) := \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{y}\|_2 \leq \epsilon\}$ . Given a set  $\mathcal{X} \subseteq \mathbb{R}^n$ , we use  $\text{cl}(\mathcal{X})$ ,  $\text{int}(\mathcal{X})$ , and  $\text{bd}(\mathcal{X})$  to denote its closure, interior, and boundary, respectively. The indicator function  $\delta_{\mathcal{X}}$  of the set  $\mathcal{X}$  is defined through  $\delta_{\mathcal{X}}(\mathbf{x}) = 0$  if  $\mathbf{x} \in \mathcal{X}$ ;  $\delta_{\mathcal{X}}(\mathbf{x}) = +\infty$  otherwise. We employ the shorthand  $\delta_{\mathbf{g}(\mathbf{x})=0}$  to compactly represent  $\delta_{\{\mathbf{x} \in \mathbb{R}^n : \mathbf{g}(\mathbf{x})=0\}}$  for any real-valued function  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ .

## 2. Assumptions and justification

In this section, we state our blanket assumptions and justify their validity in applications. To begin, we introduce the definition of the separable kernel function.

**Definition 1** *We say that  $h : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is a separable kernel function if (i)  $h$  is defined by  $h(\mathbf{x}) = \sum_{i=1}^n \varphi(x_i)$ , where  $\varphi : \mathbb{R} \rightarrow \overline{\mathbb{R}}$  is a univariate function that is continuously differentiable on  $\text{int}(\text{dom}(\varphi))$ ; (ii)  $|\varphi'(x^k)| \rightarrow +\infty$  as  $x^k \rightarrow x \in \text{bd}(\text{dom}(\varphi))$ ; (iii)  $\varphi$  is strictly convex.*

The separability structure outlined in property (i) is ubiquitous in real-world scenarios [2, 5, 20]. Properties (ii) and (iii) are referred to as Legendre-type properties, as defined in [3]. We illustrate their practicality by presenting several widely used kernel functions in Example 1.

**Assumption 1 (Problem (P))** *Let  $\text{dom}(g) = \mathcal{X}$  be a nonempty closed convex set. Suppose:*

- (i) *The function  $f$  is continuously differentiable on  $\mathcal{X}$ .*
- (ii) *The function  $g$  is convex and locally Lipschitz continuous on  $\mathcal{X}$ .*
- (iii) *There exists a strictly feasible point  $\mathbf{x}^{\text{int}} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$  and  $\mathcal{X} \subseteq \text{cl}(\text{dom}(h))$ .*
- (iv) *The function  $h$  is a separable kernel function, see Definition 1.*

Assumptions 1 (i)-(iii) or their stronger versions are widely adopted in the literature; see, e.g., Assumption A in [3], Assumption A in [5], and Definition 1 and Assumption 1 in [2]. We also refer the readers to see the practical problems that satisfy these assumptions and are solved by Bregman proximal-type methods satisfying these assumptions in, e.g., problem (3) in [23], and problem (6) in [24] and problem (7) in [12].

Next, we proceed to give assumptions of algorithm classes (i.e., the update rule (1.1)), which is characterized by the surrogate model  $\gamma$ .

**Assumption 2 (Surrogate model  $\gamma$ )** *The following hold.*

- (i) *The function  $(\mathbf{x}, \mathbf{y}) \mapsto \gamma(\mathbf{y}; \mathbf{x})$  and gradient  $(\mathbf{x}, \mathbf{y}) \mapsto \nabla \gamma(\mathbf{y}; \mathbf{x})$  (w.r.t  $\mathbf{y}$ ) are jointly continuous w.r.t.  $(\mathbf{y}, \mathbf{x})$  for all  $\mathbf{y} \in \mathcal{X}$  and  $\mathbf{x} \in \mathcal{X}$ .*

- (ii) For all  $\mathbf{x} \in \mathcal{X}$ , we have  $\nabla \gamma(\mathbf{y}; \mathbf{x})|_{\mathbf{y}=\mathbf{x}} = \nabla f(\mathbf{x})$ , and  $\gamma(\mathbf{y}; \mathbf{x})|_{\mathbf{y}=\mathbf{x}} = f(\mathbf{x})$ .
- (iii) For all  $\mathbf{x} \in \mathcal{X}$ , there exists a constant  $\bar{t} > 0$  such that  $\bar{t}\gamma(\cdot; \mathbf{x}) + h(\cdot)$  is strictly convex.
- (iv) Either  $\mathcal{X}$  is compact or we have the following condition: For all step sizes  $t \in (0, \bar{t}]$  and all sequences  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}, \{\mathbf{y}^k\}_{k \in \mathbb{N}} \subseteq \text{int}(\text{dom}(h)) \cap \mathcal{X}$  such that  $\|\mathbf{y}^k\| \rightarrow \infty$  and  $\mathbf{x}^k \rightarrow \bar{\mathbf{x}} \in \mathcal{X}$ , the following holds:

$$\lim_{k \rightarrow \infty} \gamma(\mathbf{y}^k; \mathbf{x}^k) + g(\mathbf{y}^k) + \frac{1}{t} D_h(\mathbf{y}^k, \mathbf{x}^k) = +\infty, \quad (2.1)$$

where  $D_h(\mathbf{y}, \mathbf{x}) := h(\mathbf{y}) - h(\mathbf{x}) - h'(\mathbf{x})(\mathbf{y} - \mathbf{x})$ .

Unless otherwise specified, the step size  $t$  in this paper is assumed to satisfy  $t \in (0, \bar{t}]$ .

Assumptions 2 (i) and (ii) are standard, serving to ensure the continuity and local correctness of the surrogate model  $\gamma$ . Assumption 2 (iii) usually reduces to conditions commonly adopted in the literature or is automatically satisfied for all three choices mentioned in the introduction. When  $\gamma$  represents the first-order expansion of  $f$  at the current iterate  $\mathbf{x}$ , Assumption 2 (iii) is trivially satisfied. If  $\gamma$  is selected as the original function  $f$ , it reduces to the relatively weak convexity condition, as discussed in [9, 28]. When  $\gamma$  is the second-order expansion of  $f$  at the current iterate  $\mathbf{x}$ , the  $L$ -smoothness of  $f$  and strong convexity of  $h$  suffice to guarantee Assumption 2 (iii). Assumption 2 (iv) is to ensure the well-posedness of the Bregman proximal-type update (1.1). If  $\mathcal{X}$  is unbounded, then Assumption 2 (iv) can be implied by supercoercive-type conditions, see Lemma 2 in Bauschke et al. [3] and Assumption B in Bolte et al. [9]. Interested readers are referred to Appendix B for the rigorous verification of this assumption for commonly used  $\gamma$ .

Based on Assumptions 1 and 2, we present the following lemma, which discusses the well-posedness of update (1.1) on  $\text{int}(\text{dom}(h))$ . Similar results can also be found in [3, 14].

**Lemma 2** For all  $\mathbf{x} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , we have  $T_\gamma^t(\mathbf{x}) \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , where the update mapping  $T_\gamma^t$  is defined through  $T_\gamma^t := \arg\min_{\mathbf{y}} \gamma(\mathbf{y}; \mathbf{x}^k) + g(\mathbf{y}) + \frac{1}{t} D_h(\mathbf{y}, \mathbf{x}^k)$ .

---

#### Algorithm 1: Bregman divergence-based algorithms

---

**Input:**  $\{t_k\}_{k \in \mathbb{N}}$  with  $t_k > 0$ , initial point  $\mathbf{x}^0 \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , and  $k = 0$ .

**while** Stopping Criteria is not satisfied **do**

|  $\mathbf{x}^{k+1} = T_\gamma^{t_k}(\mathbf{x}^k)$ ;  
 | Set  $k = k + 1$

**end**

**Output:** The last iterate  $\mathbf{x}^k$

---

### 3. Main results

In this section, we aim to provide a negative answer for (Q) and further establish a hardness result for the BPs. The central tool to achieve these results is a newly introduced notion, namely, *spurious stationary points*. To begin, we define the index sets

$$\mathcal{B}(\mathbf{x}) := \{b \in [n] : x_b \in \text{bd}(\text{dom}(\varphi))\}; \quad \mathcal{I}(\mathbf{x}) := \{i \in [n] : x_i \in \text{int}(\text{dom}(\varphi))\}.$$

Then, the formal definition of spurious stationary points is given as follows.

**Definition 3 (Spurious stationary points)** A point  $\mathbf{x} \in \mathcal{X}$  is defined as a spurious stationary point of problem (P) if there exists a vector  $\mathbf{p} \in \partial F(\mathbf{x})$  satisfying  $\mathbf{p}_{\mathcal{I}(\mathbf{x})} = 0$  but  $\mathbf{0} \notin \partial F(\mathbf{x})$ .

**Remark 4** (i) Spurious stationary points exist only when the kernel is non-gradient Lipschitz. (ii) For a kernel  $h$  with gradient Lipschitz property, Definition 1 (ii) implies that  $\text{dom}(h) = \mathbb{R}^n$  and  $\mathcal{I}(\mathbf{x}) = [n]$  hold for all  $\mathbf{x} \in \mathcal{X}$ , thereby precluding the existence of spurious stationary points by definition.

The following proposition demonstrates that spurious stationary points are precisely counter examples where zero does not belong to the subdifferential, yet the stationary measure equals zero. Consequently, one *cannot* classify whether the output points of the BPs are approximately stationary even when their stationarity measure is small.

**Proposition 5** Let  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \text{int}(h) \cap \mathcal{X}$  be a sequence converging to a spurious stationary point  $\tilde{\mathbf{x}}^* \in \mathcal{X}$  for problem (P). Then, the Bregman stationarity measure of  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  goes to zero, i.e., for all  $t > 0$  and  $R(\mathbf{x}) := D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})$  or  $R(\mathbf{x}) := D_h(\text{prox}_{t,\gamma}^h(\mathbf{x}), \mathbf{x})$ ,

$$\lim_{k \rightarrow \infty} R(\mathbf{x}^k) = 0.$$

**Remark 6** We consider  $R(\mathbf{x}) := D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})$  and  $R(\mathbf{x}) := D_h(\text{prox}_{t,\gamma}^h(\mathbf{x}), \mathbf{x})$  in Proposition 5 because these two measures typically evaluate the descent property of the BPs; see e.g., [3, 28]. Other measures, such as the function value gap  $F(\mathbf{x}) - F^*$  and the minimal subgradient norm  $\text{dist}(\mathbf{0}, \partial F(\mathbf{x}))$ , either fail to identify stationary points in the nonconvex setting or have no decrease guarantee for the iterate sequence of the BPs, as illustrated by Example 3.

We further uncover some practical challenges stemming from spurious stationary points. While it seems that spurious stationary points only impact the Bregman stationarity measure and that the BPs will get rid of them since they are non-stationary, the following hardness result demonstrates that we cannot escape from spurious stationary points in finite steps using Algorithm 1.

**Theorem 7 (Hardness)** If there exists a spurious stationary point  $\tilde{\mathbf{x}}^* \in \mathcal{X}$  for problem (P), then for every  $K \in \mathbb{N}$  and  $\epsilon > 0$ , there exists an initial point  $\mathbf{x}^0 \in \mathbb{B}_\epsilon(\tilde{\mathbf{x}}^*) \cap \mathcal{X} \cap \text{int}(h)$ , sufficiently close to the spurious point  $\tilde{\mathbf{x}}^*$ , such that the sequence  $\{\mathbf{x}_k\}_{k \in [K]}$  generated by Algorithm 1 satisfies

$$\mathbf{x}^k \in \mathbb{B}_\epsilon(\tilde{\mathbf{x}}^*) \quad \text{for all } k \in [K]. \quad (3.1)$$

**Remark 8** (i) If the iterates of Algorithm 1 enter a small neighborhood of a spurious stationary point, then we cannot get rid of them with any finite steps. (ii) Allowing for an infinite number of steps, under certain conditions, Algorithm 1 can eventually escape from them and converge to true stationary points. For instance, Corollary 1 in [3] establishes that, under the convexity of  $f$  and additional conditions, the sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  generated by BPG satisfies

$$f(\mathbf{x}^k) - \min_{\mathbf{x} \in \mathcal{X}} f \leq \frac{D_h(\bar{\mathbf{x}}, \mathbf{x}^0)}{t} \cdot \frac{1}{k}, \quad (3.2)$$

where  $\bar{\mathbf{x}} \in \text{argmin}_{\mathbf{x} \in \mathcal{X}} f$  is the global minimizer,  $t$  is the step size, and  $\mathbf{x}^0$  is an arbitrary initial point. The non-asymptotic convergence result (3.2) guarantees that  $f(\mathbf{x}^k) \rightarrow \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ .

Therefore, the limit points of  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  are globally optimal (stationary). (iii) Our hardness result does not contradict the non-asymptotic rate in (3.2). When (3.1) holds, the initial point  $\mathbf{x}^0$  is sufficiently close (as we constructed) to a spurious stationary point  $\tilde{\mathbf{x}}^*$ , and the distance  $D_h(\bar{\mathbf{x}}, \mathbf{x}^0)$  can be extremely large, allowing (3.2) to hold for all  $k \in [K]$ .

A remaining question is that whether spurious stationary points exist in practical optimization problems. In Appendix F, we give examples of spurious stationary points in both convex and non-convex problems. Furthermore, the following proposition demonstrates that for a broad class of convex problems with polytopal constraints, the existence of spurious points is guaranteed.

**Proposition 9 (Existence of spurious stationary points)** *Consider a convex optimization problem  $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ , where  $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  is compact,  $A \in \mathbb{R}^{m \times n}$ , and  $f$  is not constant on  $\mathcal{X}$ . If Assumption 1 holds and  $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^n$ , then every maximal point  $\tilde{\mathbf{x}}^* \in \text{argmax}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$  is a spurious stationary point.*

#### 4. Closing remark

The algorithm-dependent hardness results in this paper raise an open question:

**How can we escape spurious stationary points by modifying Bregman proximal-type algorithms and establish a valid non-asymptotic convergence rate?**

To address this question, our paper suggests one possible direction: Identify a stationarity measure that meets the equivalence (Q) and has a decreasing guarantee for the BP iterates. For convex settings, a well-known candidate stationarity measure is the function value gap  $F - F^*$ . However, when the objective functions are nonconvex, the possibility of addressing this question remains unanswered. By seeking an eligible stationarity measure, one might discover new principles for algorithm design to eliminate undesirable points.

#### References

- [1] Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory of computing*, 8(1):121–164, 2012.
- [2] Waiss Azizian, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. On the rate of convergence of Bregman proximal methods in constrained variational inequalities. *arXiv preprint arXiv:2211.08043*, 2022.
- [3] Heinz H. Bauschke, Jérôme Bolte, and Marc Teboulle. A descent lemma beyond Lipschitz gradient continuity: First-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.
- [4] Heinz H. Bauschke, Minh N. Dao, and Scott B. Lindstrom. Regularizing with Bregman–Moreau envelopes. *SIAM Journal on Optimization*, 28(4):3208–3228, 2018.
- [5] Heinz H. Bauschke, Jérôme Bolte, Jiawei Chen, Marc Teboulle, and Xianfu Wang. On linear convergence of non-Euclidean gradient methods without strong convexity and Lipschitz gradient continuity. *Journal of Optimization Theory and Applications*, 182(3):1068–1087, 2019.

- [6] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [7] Amrit Singh Bedi, Souradip Chakraborty, Anjaly Parayil, Brian M. Sadler, Pratap Tokekar, and Alec Koppel. On the hidden biases of policy mirror ascent in continuous action spaces. In *Proceedings of the 39th International Conference on Machine Learning (ICML 2022)*, pages 1716–1731. PMLR, 2022.
- [8] Benjamin Birnbaum, Nikhil R Devanur, and Lin Xiao. Distributed algorithms via gradient descent for Fisher markets. In *Proceedings of the 12th ACM conference on Electronic commerce*, pages 127–136, 2011.
- [9] Jérôme Bolte, Shoham Sabach, Marc Teboulle, and Yakov Vaisbourd. First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems. *SIAM Journal on Optimization*, 28(3):2131–2151, 2018.
- [10] Yair Censor and Stavros Andrea Zenios. Proximal minimization algorithm with D-functions. *Journal of Optimization Theory and Applications*, 73(3):451–464, 1992.
- [11] Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using Bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.
- [12] Liqun Chen, Zhe Gan, Yu Cheng, Linjie Li, Lawrence Carin, and Jingjing Liu. Graph optimal transport for cross-domain alignment. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*, pages 1542–1553. PMLR, 2020.
- [13] Nikita Doikov and Yurii Nesterov. Gradient regularization of newton method with Bregman distances. *Mathematical Programming*, 204(1):1–25, 2023.
- [14] Radu-Alexandru Dragomir, Adrien B Taylor, Alexandre d’Aspremont, and Jérôme Bolte. Optimal complexity and certification of Bregman first-order methods. *Mathematical Programming*, 194(1):41–83, 2022.
- [15] Feihu Huang, Shangqian Gao, and Heng Huang. Bregman gradient policy optimization. In *Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)*, 2022.
- [16] Feihu Huang, Junyi Li, Shangqian Gao, and Heng Huang. Enhanced bilevel optimization via Bregman distance. In *Advances in Neural Information Processing Systems 35*, pages 28928–28939, 2022.
- [17] Krzysztof C. Kiwiel. Proximal minimization methods with generalized Bregman functions. *SIAM journal on Control and Optimization*, 35(4):1142–1168, 1997.
- [18] Puya Latafat, Andreas Themelis, Masoud Ahookhosh, and Panagiotis Patrinos. Bregman Finito/MISO for nonconvex regularized finite sum minimization without Lipschitz gradient continuity. *SIAM Journal on Optimization*, 32(3):2230–2262, 2022.
- [19] Tim Tsz-Kit Lau and Han Liu. Bregman proximal Langevin Monte Carlo via Bregman-Moreau Envelopes. In *Proceedings of the 39th International Conference on Machine Learning (ICML 2022)*, pages 12049–12077. PMLR, 2022.

- [20] Jiajin Li, Jianheng Tang, Lemin Kong, Huikang Liu, Jia Li, Anthony Man-Cho So, and Jose Blanchet. A convergent single-loop algorithm for relaxation of gromov-wasserstein in graph data. In *Proceedings of the 11th International Conference on Learning Representations (ICLR 2023)*, 2023.
- [21] Wenqiang Pu, Shahana Ibrahim, Xiao Fu, and Mingyi Hong. Stochastic mirror descent for low-rank tensor decomposition under non-euclidean losses. *IEEE Transactions on Signal Processing*, 70:1803–1818, 2022.
- [22] R. Tyrrell Rockafellar and Roger J-B Wets. *Variational Analysis*, volume 317. Springer Science & Business Media, 2009.
- [23] Hongteng Xu, Dixin Luo, and Lawrence Carin. Scalable Gromov-Wasserstein learning for graph partitioning and matching. In *Advances in Neural Information Processing Systems 32*, 2019.
- [24] Hongteng Xu, Dixin Luo, Hongyuan Zha, and Lawrence Carin Duke. Gromov-Wasserstein learning for graph matching and node embedding. In *Proceedings of the 36th International Conference on Machine Learning (ICML 2019)*, pages 6932–6941. PMLR, 2019.
- [25] Long Yang, Yu Zhang, Gang Zheng, Qian Zheng, Pengfei Li, Jianhang Huang, and Gang Pan. Policy optimization with stochastic mirror descent. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence (AAAI 2022)*, volume 36, pages 8823–8831, 2022.
- [26] Hui Zhang, Yu-Hong Dai, Lei Guo, and Wei Peng. Proximal-like incremental aggregated gradient method with linear convergence under Bregman distance growth conditions. *Mathematics of Operations Research*, 46(1):61–81, 2021.
- [27] Junyu Zhang. Stochastic bergman proximal gradient method revisited: Kernel conditioning and painless variance reduction. *arXiv preprint arXiv:2401.03155*, 2024.
- [28] Siqi Zhang and Niao He. On the convergence rate of stochastic mirror descent for nonsmooth nonconvex optimization. *arXiv preprint arXiv:1806.04781*, 2018.
- [29] Daoli Zhu, Sien Deng, Minghua Li, and Lei Zhao. Level-set subdifferential error bounds and linear convergence of Bregman proximal gradient method. *Journal of Optimization Theory and Applications*, 189(3):889–918, 2021.



The appendix is structured as follows. Firstly, in Appendix A, we present some basic definitions essential for our subsequent discussions. Following that, in Appendix B, we offer an exhaustive verification of our assumptions and prove the well-posedness of Algorithm 1 under them. In Appendix C, we introduce the extended Bregman stationarity measure as a tool to investigate (Q) and unveil spurious stationary points. In Appendix D, we provide the proof of the continuity of extended stationarity measure. Then, we furnish the proofs of main results in Appendix E. Finally, we give several examples of spurious stationary points in Appendix F.

## Appendix A. Supplementary definitions

We begin by revisiting two types of subgradients for convex functions, crucial for analyzing the optimality condition.

**Definition 10 (Subgradients of convex functions)** *Rockafellar and Wets [22, Definition 8.3]* For a convex function  $l : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , we define the subgradient at the point  $\bar{x}$  by

$$\partial l(\bar{x}) := \left\{ \mathbf{v} \in \mathbb{R}^n : l(\mathbf{x}) - l(\bar{x}) \geq \mathbf{v}^\top (\mathbf{x} - \bar{x}), \text{ for all } \mathbf{x} \in \text{dom}(l) \right\},$$

and the horizon subgradient at the point  $\bar{x}$  by

$$\partial^\infty l(\bar{x}) := \left\{ \mathbf{v} \in \mathbb{R}^n : \lambda^k \mathbf{v}^k \rightarrow \mathbf{v} \text{ with } \lambda^k \rightarrow 0^+, \mathbf{v}^k \in \partial l(\mathbf{x}^k), \text{ and } \mathbf{x}^k \rightarrow \bar{x} \right\}.$$

Then, we introduce the normal cone to handle the convex constraint.

**Definition 11 (Normal cone of convex sets)** [22, Theorem 6.9] For a convex set  $\mathcal{D} \subseteq \mathbb{R}^n$ , we define the normal cone at the point  $\bar{x} \in \mathcal{D}$  via

$$\mathcal{N}_{\mathcal{D}}(\bar{x}) := \left\{ \mathbf{v} \in \mathbb{R}^n : \mathbf{v}^\top (\mathbf{x} - \bar{x}) \leq 0, \text{ for all } \mathbf{x} \in \mathcal{D} \right\}.$$

For more properties of  $\partial l$ ,  $\partial^\infty l$ , and  $\mathcal{N}_{\mathcal{D}}$ , we refer interested readers to [22, Chapter 8].

Let us revisit the three-point identity for Bregman divergence, as outlined in Lemma 3.1 of [11], i.e.,

$$D_\varphi(z, x) + D_\varphi(x, y) - D_\varphi(z, y) = (z - x) \cdot (\varphi'(y) - \varphi'(x)), \text{ for all } x, y, z \in \text{dom}(\varphi).$$

Since the strict convexity of  $\varphi$  implies the strict monotonic increase of  $\varphi'$ , the following fact holds.

**Fact 1** For a strictly convex function  $\varphi$ , the following statements hold:

- (i) If  $z \leq x \leq y$ , then  $D_\varphi(z, x) \leq D_\varphi(z, y)$ .
- (ii) If  $z \leq x \leq y$ , then  $D_\varphi(x, y) \leq D_\varphi(z, y)$ .

Next, we introduce the notion of the supercoercive property, commonly employed to ensure the well-posedness of the BPG algorithm; see [3, 9].

**Definition 12 (Supercoercive property)** Given a function  $q : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , we say that  $q$  is supercoercive if for all sequences  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  such that  $\|\mathbf{x}^k\| \rightarrow \infty$ , it holds that

$$\lim_{k \rightarrow \infty} \frac{q(\mathbf{x}^k)}{\|\mathbf{x}^k\|} \rightarrow +\infty.$$

Finally, we give the definition of Bregman proximal mapping, which serves as a crucial tool in connecting our stationarity measure with existing ones.

**Definition 13 (Bregman proximal mapping)** [4, 19] For  $t > 0$  and a kernel function  $h : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , the Bregman proximal mapping for  $F : \text{int}(\text{dom}(h)) \rightarrow \overline{\mathbb{R}}$  is defined by

$$\text{prox}_{h,F}^t(\mathbf{x}) = \underset{\mathbf{y} \in \mathbb{R}^n}{\text{argmin}} \left\{ F(\mathbf{y}) + \frac{1}{t} D_h(\mathbf{y}, \mathbf{x}) \right\}.$$

## Appendix B. Justification of Assumptions

This section is devoted to justifying our assumptions. We first present commonly used kernels that meet Definition 1 to demonstrate the generality of the separable kernels.

### Example 1

- (i) Boltzmann–Shannon entropy kernel  $h(\mathbf{x}) = \sum_{i=1}^n x_i \log(x_i)$ ;
- (ii) Fermi–Dirac entropy kernel  $h(\mathbf{x}) = \sum_{i=1}^n x_i \log(x_i) + (1 - x_i) \log(1 - x_i)$ ;
- (iii) Burg entropy kernel  $h(\mathbf{x}) = \sum_{i=1}^n -\log(x_i)$ ;
- (iv) Fractional power kernel  $h(\mathbf{x}) = \sum_{i=1}^n p x_i - \frac{x_i^p}{1-p}$  ( $0 < p < 1$ );
- (v) Hellinger entropy kernel  $h(\mathbf{x}) = \sum_{i=1}^n -\sqrt{1 - x_i^2}$ .

The implication of the separable structure is that  $\text{cl}(\text{dom}(h))$  forms a box, i.e.,

$$\text{cl}(\text{dom}(h)) = [a, c] \times [a, c] \times \cdots \times [a, c],$$

where  $a, c \in \mathbb{R} \cup \{\pm\infty\}$  with  $\text{cl}(\text{dom}(\varphi)) = [a, c]$ . Due to the convexity of  $\varphi$  and Definition 1 (ii),  $\varphi'$  is monotonically increasing and further  $\varphi'(x) \rightarrow -\infty$  (resp.  $+\infty$ ) if  $x \rightarrow a^+$  (resp.  $x \rightarrow c^-$ ) and  $a > -\infty$  (resp.  $c < +\infty$ ).

### B.1. Verification of Assumption 2 (iv)

In this subsection, we aim to identify the condition for Assumption 2 to hold.

- (i) When  $\text{dom}(\phi)$  is open, e.g.,  $\phi(x) = 1/x$ , to ensure the well-posedness of the Bregman proximal-type algorithms, we should invoke the compactness of  $\mathcal{X}$  as stated in Assumption 2 since the condition (2.1) typically fails. Such a supplement is also made in classical Bregman literature; see condition (i) in [3, Lemma 2].
- (ii) Even without  $\mathcal{X}$  being compact, the closeness of  $\text{dom}(\phi)$  and the supercoercivity of local update model is able to ensure that Assumption 2 (iv) is met, under Assumptions 1 and 2 (i)–(iii), see Proposition 14 for further details.
- (iii) Without Assumption 2 (iv), the update (1.1) may not have an optimal solution, as illustrated by Example 2.

**Example 2 (Necessity of Assumption 2 (iv))** Let  $h(x) = 1/x$ ,  $g(x) = \delta_{\mathbb{R}_+}(x)$ ,  $f(x) = (x - 3)^2$ , and  $\gamma(y; x) = (x - 3)^2 + 2(x - 3)(y - x)$ . Assumption 1 and Assumptions 2 (i)–(iii) hold for this problem with any  $\bar{t} > 0$ . We can verify that Assumption 2 (iv) fails when we pick  $x^k = \frac{1}{2} + \frac{1}{2^k}$ ,  $y^k = 2^k$  and  $t = 1$ . It turns out that the well-posedness of (1.1) does not hold. For all  $x \in [\frac{1}{2}, 1]$  and  $t = 1$ , we have

$$\begin{aligned} & \operatorname{argmin}_{y \in \mathbb{R}_+} \left\{ \gamma(y; x) + \frac{1}{t} D_h(y, x) \right\} \\ &= \operatorname{argmin}_{y \in \mathbb{R}_+} \left\{ \frac{1}{y} + \left( 2x + \frac{1}{x^2} - 6 \right) (y - x) \right\} = \emptyset, \end{aligned}$$

where the last equality is due to  $2x + \frac{1}{x^2} - 6 \leq -1$  for  $x \in [\frac{1}{2}, 1]$ .

**Proposition 14** Suppose that  $\operatorname{dom}(\phi)$  is closed. The following hold:

- (i) If the surrogate model takes the form  $\gamma(\mathbf{y}; \mathbf{x}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$  and  $h + tg$  is supercoercive for all  $t > 0$ , then Assumption 2 (iv) is satisfied.
- (ii) If the surrogate model takes the form  $\gamma(\mathbf{y}; \mathbf{x}) = f(\mathbf{y})$  and  $h + tF$  is supercoercive for all  $t > 0$ , then Assumption 2 (iv) is satisfied.
- (iii) If the surrogate model takes the form  $\gamma(\mathbf{y}; \mathbf{x}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{x}) (\mathbf{y} - \mathbf{x})$ ,  $h + tg$  is supercoercive for all  $t > 0$ , and  $f$  is a convex function, then Assumption 2 (iv) holds.

**Remark 15** Regarding (iii), the convexity condition for  $f$  is also posited in Doikov and Nesterov [13].

**Proof** Case (i):

To justify that Assumption 2 (iv) is satisfied, we proceed to prove a stronger counterpart, i.e.,

$$\lim_{k \rightarrow \infty} \frac{\gamma(\mathbf{y}^k; \mathbf{x}^k) + g(\mathbf{y}^k) + \frac{1}{t} D_h(\mathbf{y}^k, \mathbf{x}^k)}{\|\mathbf{y}^k\|} = +\infty.$$

When  $\|\mathbf{y}^k\| \rightarrow \infty$ ,  $\mathbf{x}^k \rightarrow \bar{\mathbf{x}}$ , and continuity of  $\nabla f$ , we have

$$\lim_{k \rightarrow \infty} \frac{\gamma(\mathbf{y}^k; \mathbf{x}^k)}{\|\mathbf{y}^k\|} = \nabla f(\bar{\mathbf{x}})^\top \lim_{k \rightarrow \infty} \frac{\mathbf{y}^k}{\|\mathbf{y}^k\|} < +\infty.$$

Subsequently, it remains to show

$$\lim_{k \rightarrow \infty} \frac{\frac{1}{t} D_h(\mathbf{y}^k, \mathbf{x}^k) + g(\mathbf{y}^k)}{\|\mathbf{y}^k\|} = +\infty.$$

To do so, we consider the interior coordinates and boundary coordinates separately. Recall that  $\mathcal{B}(\bar{\mathbf{x}}) = \{b \in [n] : \bar{x}_b \in \operatorname{bd}(\operatorname{dom}(\varphi))\}$  and  $\mathcal{I}(\bar{\mathbf{x}}) = \{i \in [n] : \bar{x}_i \in \operatorname{int}(\operatorname{dom}(\varphi))\}$ . Without loss of generality (WLOG), we assume that  $\operatorname{dom}(\varphi) = [a, b]$ .

- (a) For all  $i \in \mathcal{I}(\bar{\mathbf{x}})$ , we have  $\lim_{k \rightarrow \infty} |\varphi'(x_i^k)| = |\varphi'(\bar{x}_i)| < \infty$  due to the continuous differentiability of  $\varphi$ . Then, we have

$$\begin{aligned} \lim_{k \rightarrow \infty} \frac{D_\varphi(y_i^k, x_i^k)}{\|\mathbf{y}^k\|} &= \lim_{k \rightarrow \infty} \frac{\varphi(y_i^k) - \varphi(x_i^k) - \varphi'(x_i^k)(y_i^k - x_i^k)}{\|\mathbf{y}^k\|}, \\ &= \lim_{k \rightarrow \infty} \frac{\varphi(y_i^k)}{\|\mathbf{y}^k\|} - \varphi'(\bar{x}_i) \cdot \lim_{k \rightarrow \infty} \frac{y_i^k}{\|\mathbf{y}^k\|}, \quad \forall i \in \mathcal{I}(\bar{\mathbf{x}}), \end{aligned} \quad (\text{B.1})$$

where the first equality follows from the definition of  $D_\varphi$  and the second equality is due to the finiteness of  $\varphi(\bar{x}_i)$  and  $\bar{x}_i$ .

- (b) For all  $b \in \mathcal{B}(\bar{\mathbf{x}})$ , we may assume WLOG that  $\bar{x}_b = a$ . As  $\lim_{k \rightarrow \infty} x_b^k = \bar{x}_b = a$ , we have  $x_b^k < x_b^0$  for sufficiently large  $k$ . Then, we try to give a lower bound on  $D_\varphi(y_b^k, x_b^k)$  and analyze its limit for all  $b \in \mathcal{B}(\bar{\mathbf{x}})$ .

If  $y_b^k \leq x_b^0$ , then we have  $D_\varphi(y_b^k, x_b^0) \leq D_\varphi(a, x_b^0)$  due to Fact 1 (ii) and  $a < y_b^k \leq x_b^0$ . It follows that

$$D_\varphi(y_b^k, x_b^k) \geq 0 \geq D_\varphi(y_b^k, x_b^0) - D_\varphi(a, x_b^0), \quad \forall b \in \mathcal{B}(\bar{\mathbf{x}}).$$

Otherwise, we have

$$D_\varphi(y_b^k, x_b^k) \geq D_\varphi(y_b^k, x_b^0), \quad \forall b \in \mathcal{B}(\bar{\mathbf{x}}),$$

due to Fact 1 (i) and  $x_b^k < x_b^0 < y_b^k$ . Consequently, we get

$$\lim_{k \rightarrow \infty} \frac{D_\varphi(y_b^k, x_b^k)}{\|\mathbf{y}^k\|} \geq \lim_{k \rightarrow \infty} \frac{D_\varphi(y_b^k, x_b^0)}{\|\mathbf{y}^k\|} = \lim_{k \rightarrow \infty} \frac{\varphi(y_b^k) - \varphi'(x_b^0)y_b^k}{\|\mathbf{y}^k\|}, \quad \forall b \in \mathcal{B}(\bar{\mathbf{x}}), \quad (\text{B.2})$$

where the first inequality follows from the finiteness of  $D_\varphi(a, x_b^0)$ .

Together (B.1) and (B.2), we have

$$\begin{aligned} &\lim_{k \rightarrow \infty} \frac{\frac{1}{t} D_h(\mathbf{y}^k, \mathbf{x}^k) + g(\mathbf{y}^k)}{\|\mathbf{y}^k\|} \\ &\geq \lim_{k \rightarrow \infty} \frac{\frac{1}{t} h(\mathbf{y}^k) + g(\mathbf{y}^k)}{\|\mathbf{y}^k\|} - \frac{1}{t} \sum_{i \in \mathcal{I}(\bar{\mathbf{x}})} \varphi'(\bar{x}_i) \cdot \lim_{k \rightarrow \infty} \frac{y_i^k}{\|\mathbf{y}^k\|} - \frac{1}{t} \sum_{b \in \mathcal{B}(\bar{\mathbf{x}})} \varphi'(x_b^0) \cdot \lim_{k \rightarrow \infty} \frac{y_b^k}{\|\mathbf{y}^k\|} = +\infty, \end{aligned}$$

where the equality holds since  $h + tf$  is supercoercive for all  $t > 0$  and the finiteness of  $\varphi'(x_b^0)$  and  $\varphi'(\bar{x}_i)$  for all  $b \in \mathcal{B}(\bar{\mathbf{x}})$  and  $i \in \mathcal{I}(\bar{\mathbf{x}})$ . This completes the proof of case (i).

Case (ii): The ideas of the proof of case (ii) are similar to those of case (i). We just need to replace  $g$  with  $F$ .

Case (iii): Given the convexity of  $f$ , we find that  $(\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$  for any  $\mathbf{y} \in \mathbb{R}^n$ . Thus, it suffices to show

$$\lim_{k \rightarrow \infty} f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^\top (\mathbf{y}^k - \mathbf{x}^k) + g(\mathbf{y}^k) + \frac{1}{t} D_h(\mathbf{y}^k, \mathbf{x}^k) = +\infty,$$

which reduces to the case (i). We complete the proof. ■

## B.2. Proof of Lemma 2

With Assumptions 1 and 2 justified, we prove Lemma 2 to establish the well-posedness of Algorithm 1 under them.

**Proof** We prove the result by contradiction. Suppose that  $\mathbf{x} \in \text{int}(\text{dom}(h))$  and  $T_\gamma^t(\mathbf{x}) \notin \text{int}(\text{dom}(h))$ . Then, can select an interior point  $\mathbf{x}^{\text{int}} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , where its distance from  $T_\gamma^t(\mathbf{x})$  is bounded. Subsequently, we construct an intermediate point as

$$\mathbf{x}^\theta := \theta \mathbf{x}^{\text{int}} + (1 - \theta) T_\gamma^t(\mathbf{x}),$$

and an univariate function as

$$\phi(\theta) := \gamma(\mathbf{x}^\theta; \mathbf{x}) + \frac{1}{t} D_h(\mathbf{x}^\theta, \mathbf{x}) + g(\mathbf{x}^\theta),$$

where  $\theta \in (0, 1]$ .

Our strategy is to show that

$$\lim_{\theta \rightarrow 0^+} \frac{\phi(\theta) - \phi(0)}{\theta} = -\infty$$

holds<sup>1</sup>. Then, it is sufficient to demonstrate that there exists a constant  $\theta \in (0, 1)$  such that  $\phi(\theta) - \phi(0)$ , thereby contradicting the definition of  $T_\gamma^t(\mathbf{x})$ .

Now, we focus on the decomposition of  $\phi(\theta) - \phi(0)$  as below:

$$\begin{aligned} & \phi(\theta) - \phi(0) \\ &= \gamma(\mathbf{x}^\theta; \mathbf{x}) - \gamma(T_\gamma^t(\mathbf{x}); \mathbf{x}) + g(\mathbf{x}^\theta) - g(T_\gamma^t(\mathbf{x})) + \frac{1}{t} (D_h(\mathbf{x}^\theta, \mathbf{x}) - D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})) \\ &= \left[ \gamma(\mathbf{x}^\theta; \mathbf{x}) - \gamma(T_\gamma^t(\mathbf{x}); \mathbf{x}) \right] + \left[ g(\mathbf{x}^\theta) - g(T_\gamma^t(\mathbf{x})) \right] + \left[ \frac{1}{t} \sum_{j=1}^n (\varphi(x_j^\theta) - \varphi(T_\gamma^t(\mathbf{x})_j)) \right] \\ & \quad \left[ -\frac{\theta}{t} \sum_{j=1}^n \varphi'(x_j) (x_j^{\text{int}} - T_\gamma^t(\mathbf{x})_j) \right]. \end{aligned}$$

Then, we address each component individually.

(i) Due to the differentiability of  $\gamma(\cdot, \mathbf{x})$  at the point  $\mathbf{x}$  (see Assumption 2 (ii)), there exists a sufficient small  $\theta$  such that

$$\left| \frac{\gamma(\mathbf{x}^\theta; \mathbf{x}) - \gamma(T_\gamma^t(\mathbf{x}); \mathbf{x})}{\theta} \right| \leq \mathcal{O}(1).$$

(ii) Given the local Lipschitz continuity of  $g$  on  $\mathcal{X}$ , there also exists a sufficiently small  $\theta$  (i.e., you can certainly choose the same  $\theta$  to satisfy both (i) and (ii)) such that

$$\left| \frac{g(\mathbf{x}^\theta) - g(T_\gamma^t(\mathbf{x}))}{\theta} \right| \leq \mathcal{O}(1).$$

1.  $\theta \rightarrow 0^+$  denotes the limit of  $\theta$  as it approaches 0 from the right.

(iii) By Mean Value Theorem, we have

$$\sum_{j=1}^n \frac{\varphi(x_j^\theta) - \varphi(T_\gamma^t(\mathbf{x})_j)}{\theta} = \sum_{j=1}^n \varphi'(z_j^\theta) \cdot (x_j^{\text{int}} - T_\gamma^t(\mathbf{x})_j),$$

where  $z_j^\theta$  is in the interval between  $x_j^\theta$  and  $T_\gamma^t(\mathbf{x})_j$  for all  $j \in [n]$ . When  $\theta \rightarrow 0_+$ , the term  $(x_j^{\text{int}} - T_\gamma^t(\mathbf{x})_j)$  is always bounded. Moreover, we have  $z_b^\theta \rightarrow T_\gamma^t(\mathbf{x})_b$  when  $\theta \rightarrow 0^+$ . Thus, the key lies in the boundedness of  $\varphi'(T_\gamma^t(\mathbf{x})_j)$  for all  $j \in [n]$ .

We will discuss the interior coordinates and boundary coordinate of  $T_\gamma^t(\mathbf{x})$  separately.

- (a) For all  $i \in \mathcal{I}(T_\gamma^t(\mathbf{x}))$ , we know that  $T_\gamma^t(\mathbf{x})_i \in \text{int}(\text{dom}(\varphi))$ , and therefore  $\varphi'(T_\gamma^t(\mathbf{x})_i)$  is bounded due to the continuity of  $\varphi'$ .
- (b) For all  $b \in \mathcal{B}(T_\gamma^t(\mathbf{x}))$ , we have  $T_\gamma^t(\mathbf{x})_b \in \text{bd}(\text{dom}(\varphi))$  and thus  $|\varphi'(T_\gamma^t(\mathbf{x})_b)| = +\infty$  from Definition 1 (ii). WLOG, we can assume that  $\text{cl}(\text{dom}(\varphi)) = [a, c]$  and  $T_\gamma^t(\mathbf{x})_b = a$ . Moreover, as  $\varphi'$  is strictly increasing by the strict convexity of  $\varphi$ , we have  $\varphi'(T_\gamma^t(\mathbf{x})_b) = -\infty$ . Together with the fact  $x_b^{\text{int}} - T_\gamma^t(\mathbf{x})_b > 0$ , we have

$$\lim_{\theta \rightarrow 0_+} \varphi'(z_b^\theta) \cdot (x_b^{\text{int}} - T_\gamma^t(\mathbf{x})_b) = -\infty.$$

(iv) The final term is constant w.r.t.  $\theta$  and thus bounded for sure.

Putting all the pieces together, we conclude the proof.  $\blacksquare$

## Appendix C. Extended Bregman stationarity measure

To investigate (Q) and develop tools for studying spurious stationary points, we define an extended Bregman stationarity measure in this section.

To begin, we unify existing stationarity measures as  $R_\gamma^t(\mathbf{x}) := D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})$ ,<sup>2</sup> where  $t > 0$  is the step size. Conceptually, we use the relative change under the Bregman geometry to quantify the stationarity, which has been well explored in the literature [7, 15, 16, 18]. If we set  $\gamma = f$ , then the update mapping  $T_\gamma^t(\mathbf{x})$  coincides with  $\text{prox}_{h,F}^t(\mathbf{x})$ , and thus  $R_\gamma^t$  recovers the stationarity gap  $D_h(\text{prox}_{h,F}^t(\mathbf{x}), \mathbf{x})$  proposed by [28].

Despite the unification,  $R_\gamma^t$  is still not well-defined on the boundary  $\text{bd}(\text{dom}(h))$ , as the mapping  $\mathbf{x} \mapsto T_\gamma^t(\mathbf{x})$  defined by (1.1) involves the Bregman divergence function  $(\mathbf{y}, \mathbf{x}) \mapsto D_h(\mathbf{y}, \mathbf{x})$ , which is only defined on  $\text{dom}(h) \times \text{int}(\text{dom}(h))$ . Given this limitation, we are motivated to extend the domain of Bregman stationarity measures to  $\text{cl}(\text{dom}(h))$ . As a preliminary action, we define the extended update mapping that is applicable to  $\text{cl}(\text{dom}(h)) \cap \mathcal{X}$  as follows.

**Definition 16** We define the extended update mapping by  $\bar{T}_\gamma^t(\mathbf{x}) := \underset{\mathbf{y} \in \mathcal{X}}{\text{argmin}} G_\gamma^t(\mathbf{y}; \mathbf{x})$ , where

$$G_\gamma^t(\mathbf{y}; \mathbf{x}) := \gamma(\mathbf{y}; \mathbf{x}) + g(\mathbf{y}) + \underbrace{\frac{1}{t} \sum_{i \in \mathcal{I}(\mathbf{x})} D_\varphi(y_i, x_i)}_{\text{Interior coordinates}} + \underbrace{\delta_{\mathbf{y}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}}(\mathbf{y})}_{\text{Boundary coordinates}}.$$

2. Here, we ignore the difference between  $D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})$  and  $D_h(\mathbf{x}, T_\gamma^t(\mathbf{x}))$  for simplicity. As we will discuss in Remark 29, this difference would not affect our theoretical results.

This newly updated rule distinguishes between interior and boundary coordinates. For boundary coordinates of  $\mathbf{x}$ , we enforce  $\mathbf{y}$  to be equal to  $\mathbf{x}$ , while for interior coordinates, we update following the original update rule (1.1). This construction enables us to focus on the boundary points to address non-gradient Lipschitz kernel functions.

Armed with the extended update mapping  $\bar{T}_\gamma^t$ , we are ready to provide the formal definition of the extended Bregman stationarity measure  $\bar{R}_\gamma^t$ . It is worth noting that  $\bar{R}_\gamma^t$  is defined over the entire domain  $\mathcal{X}$ , and it retrieves  $R_\gamma^t(\mathbf{x})$  for  $\mathbf{x} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ .

**Definition 17 (Extended stationarity measure)** *We define the extended Bregman stationarity measure  $\bar{R}_\gamma^t(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$  as  $\bar{R}_\gamma^t(\mathbf{x}) := \sum_{i \in \mathcal{I}(\mathbf{x})} D_\varphi(\bar{T}_\gamma^t(\mathbf{x})_i, x_i)$ .*

### C.1. Well-definedness

To prove the well-definedness of the extended measure  $\bar{R}_\gamma^t$ , it suffices to prove the well-definedness of the extended mapping  $\bar{T}_\gamma^t$ . For this purpose, we first prove that the level sets of  $G_\gamma^t(\cdot; \mathbf{x})$  are nonempty; see Lemma 18. Then, due to Assumptions 2 (i) and (iv), we see that the level sets of  $G_\gamma^t(\cdot; \mathbf{x})$  are compact. Further, Assumption 2 (iii) guarantees the strict convexity of  $G_\gamma^t(\cdot; \mathbf{x})$ . Both the strict convexity and level boundedness of  $G_\gamma^t$  ensure the well-definedness and uniqueness of the extended update mapping  $\bar{T}_\gamma^t$ .

**Lemma 18 (Existence)** *For all  $\mathbf{x} \in \mathcal{X}$ , we have*

$$\text{argmin}_{\mathbf{y} \in \mathcal{X}} G_\gamma^t(\mathbf{y}; \mathbf{x}) \neq \emptyset.$$

**Proof** To ensure the existence of the optimal solution, it suffices to show that, for some constant  $c \in \mathbb{R}$ , the level set  $\{\mathbf{y} \in \mathcal{X} : G_\gamma^t(\mathbf{y}; \mathbf{x}) \leq c\}$  is compact for any  $\mathbf{x} \in \mathcal{X}$ . At first, from [22, Theorem 1.6], we know this level set is closed due to the lower semicontinuity of  $G_\gamma^t(\cdot; \mathbf{x})$  for any  $\mathbf{x} \in \mathcal{X}$ .

Next, our task is to show the set  $\{\mathbf{y} \in \mathcal{X} : G_\gamma^t(\mathbf{y}; \mathbf{x}) \leq c\}$  is bounded. We consider two cases separately:

(i) If  $\mathbf{x} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , the boundedness of the set  $\{\mathbf{y} \in \mathcal{X} : G_\gamma^t(\mathbf{y}; \mathbf{x}) \leq c\}$  is from Assumption 2 (iv), as the coerciveness can imply one non-empty bounded level set.

(ii) If  $\mathbf{x} \in \text{bd}(\text{dom}(h)) \cap \mathcal{X}$ , it suffices to show the coerciveness of  $G_\gamma^t(\cdot; \mathbf{x})$ . That is, if we consider an arbitrary sequence  $\{\mathbf{y}^k\} \subseteq \mathcal{X}$  with  $\|\mathbf{y}^k\| \rightarrow +\infty$ , we have  $G_\gamma^t(\mathbf{y}^k; \mathbf{x}) \rightarrow +\infty$ .

WLOG, we can assume that  $\mathbf{y}_{\mathcal{B}(\mathbf{x})}^k \equiv \mathbf{x}_{\mathcal{B}(\mathbf{x})}$ ; otherwise,  $G_\gamma^t(\mathbf{y}^k; \mathbf{x}) = +\infty$ , for any  $k \in \mathbb{N}$ . Then, we pick up an interior point  $\mathbf{x}^{\text{int}} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$  and construct three sequences  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ ,  $\{\tilde{\mathbf{y}}^k\}_{k \in \mathbb{N}}$  and  $\{\tilde{\mathbf{x}}^k\}_{k \in \mathbb{N}}$  satisfying

$$\begin{aligned} \mathbf{x}^k &\rightarrow \mathbf{x} \text{ with } \mathcal{B}(\mathbf{x}^k) = \mathcal{B}(\mathbf{x}), \forall k \in \mathbb{N} \\ \tilde{\mathbf{y}}^k &= (1 - \theta_k)\mathbf{y}^k + \theta_k\mathbf{x}^{\text{int}}, \forall k \in \mathbb{N}, \\ \tilde{\mathbf{x}}^k &= (1 - \theta_k)\mathbf{x}^k + \theta_k\mathbf{x}^{\text{int}}, \forall k \in \mathbb{N}, \end{aligned}$$

where  $\theta_k \in (0, 1)$  and  $\theta_k \rightarrow 0$ . Then, it is trivial to conclude that  $\|\tilde{\mathbf{y}}^k\| \rightarrow +\infty$  and  $\tilde{\mathbf{x}}^k \rightarrow \mathbf{x}$ . Moreover, due to  $\tilde{\mathbf{y}}_{\mathcal{B}(\mathbf{x})}^k = \tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})}^k, \forall k \in \mathbb{N}$ , we get

$$\gamma(\tilde{\mathbf{y}}^k; \tilde{\mathbf{x}}^k) + g(\tilde{\mathbf{y}}^k) + \frac{1}{t} D_h(\tilde{\mathbf{y}}^k, \tilde{\mathbf{x}}^k) = \gamma(\tilde{\mathbf{y}}^k; \tilde{\mathbf{x}}^k) + g(\tilde{\mathbf{y}}^k) + \frac{1}{t} \sum_{i \in \mathcal{I}(\mathbf{x})} D_\varphi(\tilde{\mathbf{y}}_i^k, \tilde{\mathbf{x}}_i^k).$$

By the continuity of the mappings  $\gamma(\cdot), g(\cdot)$ , and  $D_\varphi(\cdot, \cdot)$ , there exists a sufficiently small  $\theta_k > 0$  such that

$$\left| \gamma(\tilde{\mathbf{y}}^k; \tilde{\mathbf{x}}^k) + g(\tilde{\mathbf{y}}^k) + \frac{1}{t} D_h(\tilde{\mathbf{y}}^k, \tilde{\mathbf{x}}^k) - G_\gamma^t(\mathbf{y}^k; \mathbf{x}) \right| \leq 1.$$

From Assumption 2 (iv), we know  $\gamma(\tilde{\mathbf{y}}^k; \tilde{\mathbf{x}}^k) + g(\tilde{\mathbf{y}}^k) + \frac{1}{t} D_h(\tilde{\mathbf{y}}^k, \tilde{\mathbf{x}}^k) \rightarrow +\infty$ . Thus, we can conclude that  $G_\gamma^t(\mathbf{y}^k; \mathbf{x}) \rightarrow +\infty$ . ■

## C.2. Continuity of extended stationarity measure

To investigate the equivalence described in (Q), our first step is to establish the continuity of  $\bar{R}_\gamma^t$ . This task involves recognizing that  $\bar{R}_\gamma^t$ 's definition depends on both  $\bar{T}_\gamma^t$  and  $\mathcal{I}(\mathbf{x})$ , which may exhibit discontinuity. To start, we establish the continuity of the extended update mapping:

**Proposition 19** *The extended update mapping  $\bar{T}_\gamma^t : \mathcal{X} \rightarrow \mathbb{R}^n$  is continuous on the domain  $\mathcal{X}$ .*

The continuity of  $\bar{T}_\gamma^t$  on  $\mathcal{X}$  serves as a fundamental property of the extended update mapping. Not only does it provide insight into  $\bar{T}_\gamma^t(\mathbf{x})$  for  $\mathbf{x} \in \text{bd}(\text{dom}(h)) \cap \mathcal{X}$ , but it also plays a crucial role in establishing the continuity of  $\bar{R}_\gamma^t$ . Leveraging the continuity of  $\bar{T}_\gamma^t$  and the structure of  $\mathcal{I}(\mathbf{x})$ , we establish one of the main theoretical results as below:

**Theorem 20** *The extended stationarity measure  $\bar{R}_\gamma^t : \mathcal{X} \rightarrow \mathbb{R}$  is continuous on the domain  $\mathcal{X}$ .*

## C.3. Necessity of zero extended stationarity measure

In this subsection, we will establish the necessary conditions by utilizing the extended stationarity measure. The key idea for proving necessity is through the fixed-point equation  $\bar{T}_\gamma^t(\mathbf{x}) = \mathbf{x}$ . Below is the flowchart outlining the proofs:

$$\boxed{\mathbf{0} \in \partial F(\mathbf{x})} \xrightarrow{\text{Proposition 22}} \boxed{\bar{T}_\gamma^t(\mathbf{x}) = \mathbf{x}} \xleftrightarrow{\text{Proposition 21}} \boxed{\bar{R}_\gamma^t(\mathbf{x}) = 0} \text{ for } \mathbf{x} \in \mathcal{X}.$$

We start with establishing the equivalence between  $\bar{T}_\gamma^t(\mathbf{x}) = \mathbf{x}$  and  $\bar{R}_\gamma^t(\mathbf{x}) = 0$ .

**Proposition 21** *For all  $\mathbf{x} \in \mathcal{X}$ , the extended stationarity measure being zero, i.e.,  $\bar{R}_\gamma^t(\mathbf{x}) = 0$ , is equivalent to  $\bar{T}_\gamma^t(\mathbf{x}) = \mathbf{x}$ .*



**Proof** Considering the definition of  $\overline{T}_\gamma^t(\mathbf{x})$ , we observe that  $\overline{T}_\gamma^t(\mathbf{x})_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}$  on the boundary coordinates. Regarding those interior coordinates, the definition of  $\overline{R}_\gamma^t$  guarantees that  $\overline{R}_\gamma^t(\mathbf{x}) = 0$  if and only if  $\overline{T}_\gamma^t(\mathbf{x})_{\mathcal{I}(\mathbf{x})} = \mathbf{x}_{\mathcal{I}(\mathbf{x})}$ . By combining these two observations and the fact that  $\mathcal{I}(\mathbf{x}) \cup \mathcal{B}(\mathbf{x}) = [n]$ , we establish the equivalence between  $\overline{R}_\gamma^t(\mathbf{x}) = 0$  and  $\overline{T}_\gamma^t(\mathbf{x}) = \mathbf{x}$ . ■

Next, we will investigate the relation between the fixed-point equation  $\overline{T}_\gamma^t(\mathbf{x}) = \mathbf{x}$  and the stationary condition  $\mathbf{0} \in \partial F(\mathbf{x})$ .

**Proposition 22** *If  $\mathbf{x} \in \mathcal{X}$  is a stationary point, then we have  $\overline{T}_\gamma^t(\mathbf{x}) = \mathbf{x}$ .*

**Proof** To start, it is sufficient to demonstrate that  $\mathbf{x} = \operatorname{argmin}_{\mathbf{y}} G_\gamma^t(\mathbf{y}; \mathbf{x})$ , a condition that is equivalent to  $\mathbf{0} \in \partial G_\gamma^t(\mathbf{x}; \mathbf{x})$  according to optimality condition. Based on Assumption 2 (ii) and [22, Corollary 10.9], we have

$$\partial G_\gamma^t(\mathbf{y}; \mathbf{x}) \big|_{\mathbf{y}=\mathbf{x}} = \nabla f(\mathbf{x}) + \partial g(\mathbf{x}) + \partial \delta_{\mathbf{y}_{\mathcal{B}(\mathbf{x})}=\mathbf{x}_{\mathcal{B}(\mathbf{x})}}(\mathbf{y}) \big|_{\mathbf{y}=\mathbf{x}}. \quad (\text{C.1})$$

Since we have  $\mathbf{0} \in \partial \delta_{\mathbf{y}_{\mathcal{B}(\mathbf{x})}=\mathbf{x}_{\mathcal{B}(\mathbf{x})}}(\mathbf{y}) \big|_{\mathbf{y}=\mathbf{x}}$ , we know  $\mathbf{0} \in \partial F(\mathbf{x})$  and then  $\mathbf{0} \in \partial G_\gamma^t(\mathbf{x}; \mathbf{x})$ . ■

Equipped with Proposition 21 and Proposition 22, we are poised to demonstrate our primary discovery: The extended stationarity measure equals zero for all stationary points.

**Theorem 23** *If  $\mathbf{x} \in \mathcal{X}$  is a stationary point, i.e.,  $\mathbf{0} \in \partial F(\mathbf{x})$ , then we have  $\overline{R}_\gamma^t(\mathbf{x}) = 0$ .*

Combining Theorem 23 with the continuity of  $\overline{R}_\gamma^t$  (as demonstrated in Theorem 20), we can now establish the necessity direction of equivalence (Q).

**Corollary 24** *Let the sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \operatorname{int}(\operatorname{dom}(h)) \cap \mathcal{X}$  converge to  $\overline{\mathbf{x}} \in \mathcal{X}$  with  $\mathbf{0} \in \partial F(\overline{\mathbf{x}})$ . Then, we have*

$$\lim_{k \rightarrow \infty} \overline{R}_\gamma^t(\mathbf{x}^k) = 0. \quad (\text{C.2})$$

**Proof** Here is a one-line proof:

$$\lim_{k \rightarrow \infty} \overline{R}_\gamma^t(\mathbf{x}^k) = \overline{R}_\gamma^t \left( \lim_{k \rightarrow \infty} \mathbf{x}^k \right) = \overline{R}_\gamma^t(\overline{\mathbf{x}}) = 0,$$

where the first equality follows from the continuity of  $\overline{R}_\gamma^t$  (see Theorem 20), and the last one is due to the necessity of  $\overline{R}_\gamma^t$  (see Theorem 23). ■

Corollary 24 claims that if a limit point of a sequence is stationary, then (C.2) holds true. However, existing counterparts of (C.2) are typically established under the assumption that the kernel functions are gradient Lipschitz, see [21, 25, 28]. In contrast to existing works, Corollary 24 presented in this paper can be applied to a broader class of kernel functions, such as those without the gradient Lipschitz property.

#### C.4. Non-sufficiency of zero extended stationarity measure

Based on the extended stationarity measure, we give a characterization of spurious stationary points, demonstrating the non-sufficiency of zero extended stationarity measure. Then, the proof of Proposition 5 directly follows.

**Proposition 25 (Characterization of spurious stationary points)** *A point  $\mathbf{x} \in \mathcal{X}$  is a spurious stationary point if and only if*

$$\overline{R}_\gamma^t(\mathbf{x}) = 0 \text{ but } \mathbf{0} \notin \partial F(\mathbf{x}).$$

**Proof** From Definition 3, it is sufficient to show the equivalence between  $\overline{R}_\gamma^t(\mathbf{x}) = 0$  and the existence of a vector  $\mathbf{p} \in \partial F(\mathbf{x})$  where  $\mathbf{p}_{\mathcal{I}(\mathbf{x})} = \mathbf{0}$ . Following the proof of Proposition 22, we will proceed to check the equivalent optimality condition of  $\mathbf{0} \in \partial G_\gamma^t(\mathbf{x}; \mathbf{x})$ . Observing that  $\partial \delta_{\mathbf{y}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}}(\mathbf{x}) = \text{span}\{\mathbf{e}_b : b \in \mathcal{B}(\mathbf{x})\}$ , we know that  $\mathbf{0} \in \partial G_\gamma^t(\mathbf{x}; \mathbf{x})$  is equivalent to the existence of a vector  $\mathbf{p} \in \partial F(\mathbf{x})$  with  $\mathbf{p}_{\mathcal{I}(\mathbf{x})} = \mathbf{0}$ . We complete the proof.  $\blacksquare$

Considering the continuity of  $\overline{R}_\gamma^t$ , if a sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \mathcal{X} \cap \text{int}(\text{dom}(h))$  converges to a spurious stationary point, then  $\overline{R}_\gamma^t(\mathbf{x}^k) \rightarrow 0$ . This proves Proposition 5.

#### Appendix D. Proof of continuity of extended stationarity measure

To begin, we introduce two technical lemmas essential for proving Proposition 19.

**Lemma 26** *Suppose the sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \text{int}(\text{dom}(h)) \cap \mathcal{X}$  is bounded, and the sequence  $\{T_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  converges to  $\mathbf{x}^* \in \mathcal{X}$ . We define a sequence  $\{\mathbf{d}^{k+1}\}_{k \in \mathbb{N}}$  satisfying*

$$\nabla_\gamma \left( T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k \right) + \mathbf{d}^{k+1} + \frac{1}{t} \left( \nabla h \left( T_\gamma^t(\mathbf{x}^k) \right) - \nabla h(\mathbf{x}^k) \right) = \mathbf{0} \text{ and } \mathbf{d}^{k+1} \in \partial g \left( T_\gamma^t(\mathbf{x}^k) \right). \quad (\text{D.1})$$

*Then, the sequence  $\{\mathbf{d}^{k+1}\}_{k \in \mathbb{N}}$  is bounded.*

**Proof** We prove the boundedness of  $\{\mathbf{d}^{k+1}\}_{k \in \mathbb{N}}$  by contradiction. Suppose that the sequence  $\{\mathbf{d}^{k+1}\}_{k \in \mathbb{N}}$  is not bounded. Then, there must exist a subsequence that diverges. WLOG, we can assume that  $\|\mathbf{d}^{k+1}\| \rightarrow \infty$  and  $\mathbf{d}^{k+1}/\|\mathbf{d}^{k+1}\| \rightarrow \mathbf{d}^*$  for some  $\mathbf{d}^*$ . Due to Definition 10 and  $T_\gamma^t(\mathbf{x}^k) \rightarrow \mathbf{x}^*$ , we have  $\mathbf{d}^* \in \partial^\infty g(\mathbf{x}^*)$ . Moreover, owing to the convexity and continuity of the function  $g$ , we can apply Rockafellar and Wets [22, Proposition 8.12] to get  $\partial^\infty g(\mathbf{x}^*) = \mathcal{N}_{\mathcal{X}}(\mathbf{x}^*)$  and hence we have

$$(\mathbf{x}^* - \mathbf{x})^\top \mathbf{d}^* \geq 0 \text{ for all } \mathbf{x} \in \mathcal{X}. \quad (\text{D.2})$$

In view of (D.2), we consider the following two scenarios separately:

- (i)  $(\mathbf{x}^* - \mathbf{x})^\top \mathbf{d}^* = 0$  for all  $\mathbf{x} \in \mathcal{X}$ ;
- (ii) There exists some  $\mathbf{x} \in \mathcal{X}$  such that  $(\mathbf{x}^* - \mathbf{x})^\top \mathbf{d}^* > 0$ .

**Scenarios (i):** For all  $k \in \mathbb{N}$ , we have

$$(T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^k)^\top \mathbf{d}^* = 0, \quad (\text{D.3})$$

by substituting  $\mathbf{x} = \mathbf{x}^k$  and  $\mathbf{x} = T_\gamma^t(\mathbf{x}^k)$  in (D.2), and summing them.

From (D.3), we have

$$\varphi' \left( T_\gamma^t(\mathbf{x}^k)_i \right) - \varphi'(x_i^k) = -t \left( \nabla_i \gamma(T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k) + d_i^{k+1} \right), \forall i \in [n]. \quad (\text{D.4})$$

Next, we want to argue that the left-hand side of (D.4) will go to infinity, which will contradict with (D.3).

Given the boundedness of the sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  and  $\{T_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}} \rightarrow \mathbf{x}^*$ , it follows from the joint continuity property of  $\nabla \gamma(\cdot; \cdot)$ , as stated in Assumption 2 (i), that  $\{\nabla_i \gamma(T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k)\}_{k \in \mathbb{N}}$  is also bounded. We proceed to discuss the term  $d_i^{k+1}$  for all  $i \in [n]$ . Let  $\mathcal{I}^+ := \{i \in [n] : d_i^* > 0\}$  and  $\mathcal{I}^- := \{i \in [n] : d_i^* < 0\}$ . Thus, we have  $d_i^{k+1} \rightarrow +\infty$  (resp.  $d_i^{k+1} \rightarrow -\infty$ ) by  $d_i^{k+1} / \|\mathbf{d}^{k+1}\| \rightarrow d_i^*$ , for all  $i \in \mathcal{I}^+$  (resp.  $\mathcal{I}^-$ ). Hence, the equation (D.4) yields

$$\varphi' \left( T_\gamma^t(\mathbf{x}^k)_i \right) - \varphi'(x_i^k) \rightarrow -\infty \text{ (resp. } +\infty) \text{ for } i \in \mathcal{I}^+ \text{ (resp. } i \in \mathcal{I}^-).$$

Moreover, since  $\varphi'$  is strictly increasing by the strict convexity of  $\varphi$ , for sufficiently large  $k$ , we have

$$T_\gamma^t(\mathbf{x}^k)_i < x_i^k \text{ for } i \in \mathcal{I}^+ \text{ and } T_\gamma^t(\mathbf{x}^k)_i > x_i^k \text{ for } i \in \mathcal{I}^-.$$

Consequently, we get

$$(T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^k)^\top \mathbf{d}^* < 0 \text{ for sufficiently large } k,$$

which contradicts (D.3).

**Scenarios (ii):** There must exist a point  $x \in \mathcal{X}$  and a positive constant  $\alpha > 0$  such that  $(\mathbf{x}^* - \mathbf{x})^\top \mathbf{d}^* \geq \alpha$ . Additionally, we can select some  $\mathbf{x}^{\text{int}} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , which is sufficiently close to  $\mathbf{x}$  such that  $(\mathbf{x}^{\text{int}} - \mathbf{x})^\top \mathbf{d}^* \leq \frac{\alpha}{2}$ . Therefore, we have

$$(\mathbf{x}^* - \mathbf{x}^{\text{int}})^\top \mathbf{d}^* = \lim_{k \rightarrow \infty} (T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \frac{\mathbf{d}^{k+1}}{\|\mathbf{d}^{k+1}\|} \geq \frac{\alpha}{2}.$$

where the first equality follows from  $T_\gamma^t(\mathbf{x}^k) \rightarrow \mathbf{x}^*$  and the definition of  $\mathbf{d}^*$ . Then, for sufficiently large  $k$ , we have

$$(T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \mathbf{d}^{k+1} \geq \frac{\alpha}{2} \|\mathbf{d}^{k+1}\|. \quad (\text{D.5})$$

Furthermore, by multiplying (D.1) by  $T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}}$ , we obtain

$$(T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \left[ \nabla \gamma \left( T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k \right) + \mathbf{d}^{k+1} + \frac{1}{t} \left( \nabla h(T_\gamma^t(\mathbf{x}^k)) - \nabla h(\mathbf{x}^k) \right) \right] = 0.$$

It follows that

$$\begin{aligned} & (T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \left( \nabla h(T_\gamma^t(\mathbf{x}^k)) - \nabla h(\mathbf{x}^k) \right) \\ &= - (T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \mathbf{d}^{k+1} - (T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \nabla \gamma \left( T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k \right) \\ &\leq - \frac{\alpha}{2} \|\mathbf{d}^{k+1}\|^2 - (T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top \nabla \gamma \left( T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k \right) \rightarrow -\infty, \end{aligned} \quad (\text{D.6})$$

where the first inequality is due to (D.5), and the last limit follows from the boundedness of  $\{\nabla \gamma(T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k)\}_{k \in \mathbb{N}}$  and  $\|\mathbf{d}^{k+1}\| \rightarrow \infty$ .

We now utilize (D.6) to establish the contradiction. The key proof idea follows from the one in Lemma 2. For all  $k \in \mathbb{N}$ , we construct a sequence of intermediate points as

$$\mathbf{x}^{\theta,k} := \theta \mathbf{x}^{\text{int}} + (1 - \theta) T_\gamma^t(\mathbf{x}^k),$$

and a sequence of univariate functions as

$$\phi_k(\theta) := \gamma(\mathbf{x}^{\theta,k}; \mathbf{x}) + \frac{1}{t} D_h(\mathbf{x}^{\theta,k}, \mathbf{x}) + g(\mathbf{x}^{\theta,k}),$$

where  $\theta \in (0, 1]$ . By the definition of  $T_\gamma^t(\mathbf{x}^k)$ , we have  $\phi_k(0) = \min_{\theta \in [0,1]} \phi_k(\theta)$  for  $k \in \mathbb{N}$ . Then, our strategy is to show that for sufficiently large  $k$

$$\lim_{\theta \rightarrow 0^+} \frac{\phi_k(\theta) - \phi_k(0)}{\theta} < 0,$$

and thus yields a contradiction. By Mean Value Theorem, we have

$$\begin{aligned} & \frac{\phi_k(\theta) - \phi_k(0)}{\theta} \\ &= \frac{\gamma(\mathbf{x}^{\theta,k}; \mathbf{x}) - \gamma(T_\gamma^t(\mathbf{x}^{\theta,k}); \mathbf{x})}{\theta} + \frac{1}{t} (T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^{\text{int}})^\top (\nabla h(\mathbf{z}^{\theta,k}) - \nabla h(\mathbf{x}^k)) \\ & \quad + \frac{g(\mathbf{x}^{\theta,k}) - g(T_\gamma^t(\mathbf{x}^k))}{\theta}, \end{aligned} \tag{D.7}$$

where  $\mathbf{z}^{\theta,k}$  is between  $\mathbf{x}^{\theta,k}$  and  $T_\gamma^t(\mathbf{x}^k)$ . Owing to the continuous differentiability of  $\gamma$  and local Lipschitz continuity of  $g$ , we know the first and third terms in (D.7) are bounded. From (D.6), as  $\theta$  approaches  $0_+$ , the second term tends toward  $-\infty$  for sufficiently large  $k$ . Consequently, we have

$$\lim_{k \rightarrow \infty} \lim_{\theta \rightarrow 0^+} \frac{\phi_k(\theta) - \phi_k(0)}{\theta} = -\infty,$$

which leads to a contradiction. We complete our proof.  $\blacksquare$

Next, we extend Lemma 26 to a more general case, whose proof techniques are essentially the same as the one developed in Lemma 26.

**Corollary 27** *Suppose that the sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \mathcal{X}$  is bounded satisfying  $\mathcal{I}(\mathbf{x}^k) \equiv \mathcal{I}_0 \subseteq [n]$  and  $\mathcal{B}(\mathbf{x}^k) \equiv \mathcal{B}_0$ , and the sequence  $\{\bar{T}_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  converges to  $\mathbf{x}^* \in \mathcal{X}$ . We define a sequence  $\{\mathbf{d}^{k+1}\}_{k \in \mathbb{N}}$  satisfying*

$$\nabla \gamma(\bar{T}_\gamma^t(\mathbf{x}^k); \mathbf{x}^k) + \mathbf{d}^{k+1} + \mathbf{p}^{k+1} + \frac{1}{t} \sum_{i \in \mathcal{I}_0} \nabla \left( D_\varphi(\bar{T}_\gamma^t(\mathbf{x}^k)_i, x_i^k) \right) = \mathbf{0}, \tag{D.8}$$

where  $\mathbf{d}^{k+1} \in \partial g(\bar{T}_\gamma^t(\mathbf{x}^k))$  and  $\mathbf{p}^{k+1} \in \partial \delta_{\mathbf{y}_{\mathcal{B}_0} = \mathbf{x}_{\mathcal{B}_0}^k}(\bar{T}_\gamma^t(\mathbf{x}^k))$ , for all  $k \in \mathbb{N}$ . Then, the sequence  $\{\mathbf{d}_{\mathcal{I}_0}^{k+1}\}_{k \in \mathbb{N}}$  is bounded.

**Proof** WLOG, we assume that  $\text{cl}(\text{dom}(\varphi)) = [a, c]$  and  $\mathbf{x}_{\mathcal{B}_0}^k \equiv \mathbf{x}_{\mathcal{B}_0}^0$  for all  $k \in \mathbb{N}$ . Let  $\tilde{g} = g + \delta_{\mathbf{y}_{\mathcal{B}_0} = \mathbf{x}_{\mathcal{B}_0}^0}$ ,  $\tilde{\mathcal{X}} = \mathcal{X} \cap \{\mathbf{y} \in \mathcal{X} : \mathbf{y}_{\mathcal{B}_0} = \mathbf{x}_{\mathcal{B}_0}^0\}$ , and  $\tilde{\mathbf{d}}^{k+1} = \mathbf{d}^{k+1} + \mathbf{p}^{k+1}$ . Then, we have

$$\tilde{\mathbf{d}}_{\mathcal{B}_0}^{k+1} = \mathbf{d}_{\mathcal{B}_0}^{k+1}, \quad \tilde{\mathbf{d}}^{k+1} \in \partial \tilde{g} \left( \bar{T}_\gamma^t(\mathbf{x}^k) \right), \quad \text{and} \quad \text{dom}(\tilde{g}) = \tilde{\mathcal{X}}.$$

We continue to rewrite (D.8) as

$$\nabla \gamma(\bar{T}_\gamma^t(\mathbf{x}^k); \mathbf{x}^k) + \tilde{\mathbf{d}}^{k+1} + \frac{1}{t} \sum_{i \in \mathcal{I}_0} \nabla \left( D_\varphi(\bar{T}_\gamma^t(\mathbf{x}^k)_i, x_i^k) \right) = \mathbf{0}, \quad (\text{D.9})$$

where  $\tilde{\mathbf{d}}^{k+1} \in \partial \tilde{g} \left( \bar{T}_\gamma^t(\mathbf{x}^k) \right)$ . Due to the boundedness of  $\{\bar{T}_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  and  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ , and the continuity of  $\nabla \gamma$ , it follows that the sequence  $\{\nabla \gamma(\bar{T}_\gamma^t(\mathbf{x}^k); \mathbf{x}^k)\}_{k \in \mathbb{N}}$  is bounded. Moreover, we know the sequence  $\{\tilde{\mathbf{d}}_{\mathcal{B}_0}^{k+1}\}_{k \in \mathbb{N}}$  is bounded from the boundedness of  $\{\bar{T}_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  and (D.9).

We left to show the sequence  $\{\mathbf{d}_{\mathcal{I}_0}^{k+1}\}_{k \in \mathbb{N}}$  is bounded. We prove this result by contradiction. Suppose that the sequence  $\{\mathbf{d}_{\mathcal{I}_0}^{k+1}\}_{k \in \mathbb{N}}$  is unbounded. Hence, the sequence  $\{\tilde{\mathbf{d}}^{k+1}\}_{k \in \mathbb{N}}$  is also unbounded. Then, there must exist a subsequence that diverges. WLOG, we can assume that  $\tilde{\mathbf{d}}^{k+1} / \|\tilde{\mathbf{d}}^{k+1}\| \rightarrow \bar{\mathbf{d}}$  for some  $\bar{\mathbf{d}}$ . Here, we have  $\|\bar{\mathbf{d}}_{\mathcal{I}_0}\| = 1$  and  $\bar{\mathbf{d}}_{\mathcal{B}_0} = \mathbf{0}$  due to boundedness of  $\{\tilde{\mathbf{d}}_{\mathcal{B}_0}^{k+1}\}_{k \in \mathbb{N}}$ . Then, the left proof is the same as the one in Lemma 26. We omit the proof details. ■

**Lemma 28** *If the sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \text{int}(\text{dom}(h)) \cap \mathcal{X}$  is bounded, then the sequence  $\{T_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  is also bounded.*

**Proof** We prove this result by contradiction. WLOG, we can assume  $\mathbf{x}^k \rightarrow \bar{\mathbf{x}} \in \mathcal{X}$  by the boundedness of  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ . From Assumption 2 (iv), we have

$$\lim_{k \rightarrow +\infty} \gamma(T_\gamma^t(\mathbf{x}^k), \mathbf{x}^k) + \frac{1}{t} D_h(T_\gamma^t(\mathbf{x}^k), \mathbf{x}^k) + g(T_\gamma^t(\mathbf{x}^k)) = +\infty. \quad (\text{D.10})$$

Moreover, we have

$$\lim_{k \rightarrow +\infty} \gamma(\mathbf{x}^k; \mathbf{x}^k) + \frac{1}{t} D_h(\mathbf{x}^k, \mathbf{x}^k) + g(\mathbf{x}^k) = \gamma(\bar{\mathbf{x}}; \bar{\mathbf{x}}) + g(\bar{\mathbf{x}}) = f(\bar{\mathbf{x}}) + g(\bar{\mathbf{x}}) < +\infty,$$

where the first equality is owing to the continuity of  $\gamma$  and  $g$ , the second one is from Assumption 2 (ii), and the last inequality is due to  $\bar{\mathbf{x}} \in \mathcal{X}$ . Then, together with (D.10), for  $k$  large enough, it follows that

$$\gamma(T_\gamma^t(\mathbf{x}^k), \mathbf{x}^k) + \frac{1}{t} D_h(T_\gamma^t(\mathbf{x}^k), \mathbf{x}^k) + g(T_\gamma^t(\mathbf{x}^k)) > \gamma(\mathbf{x}^k, \mathbf{x}^k) + \frac{1}{t} D_h(\mathbf{x}^k, \mathbf{x}^k) + g(\mathbf{x}^k),$$

which contradicts with the definition of  $T_\gamma^t(\mathbf{x}^k)$ . Hence, the sequence  $\{T_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  is bounded. ■

Now, we are ready to give a full proof of Proposition 19.

### D.1. Proof of Proposition 19

**Proof** To show that the mapping  $\bar{T}_\gamma^t(\cdot)$  is continuous, it suffices to show that for any sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \mathcal{X}$  converging to  $\mathbf{x} \in \mathcal{X}$ , it holds that:

$$\lim_{k \rightarrow \infty} \bar{T}_\gamma^t(\mathbf{x}^k) = \bar{T}_\gamma^t(\mathbf{x}).$$

To proceed, we consider the following two scenarios sequentially:

- (i) The sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \text{int}(\text{dom}(h)) \cap \mathcal{X}$ .
- (ii) The sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}} \subseteq \text{cl}(\text{dom}(h)) \cap \mathcal{X}$ .

To start with, we consider the simple case (i). Later on, we can extend our proof strategy to consider the general case by considering the interior and boundary coordinates of  $\mathbf{x}$  separately. Then, we can reduce the general case to the simple case considered here.

**Scenarios (i):** Due to  $\bar{T}_\gamma^t(\cdot) = T_\gamma^t(\cdot)$  on  $\text{int}(\text{dom}(h)) \cap \mathcal{X}$ . It is equivalent to show

$$\bar{T}_\gamma^t(\mathbf{x}) = \lim_{k \rightarrow \infty} T_\gamma^t(\mathbf{x}^k).$$

The remainder of the proof proceeds in two steps. We give a proof sketch here initially.

- **Step 1:** For the boundary coordinates of  $\mathbf{x}$ , we have  $\lim_{k \rightarrow \infty} T_\gamma^t(\mathbf{x}^k)_{\mathcal{B}(\mathbf{x})} = \bar{T}_\gamma^t(\mathbf{x})_{\mathcal{B}(\mathbf{x})}$ .

By Lemma 28, we can pass to a subsequence such that  $T_\gamma^t(\mathbf{x}^k) \rightarrow \tilde{\mathbf{x}} \in \text{cl}(\text{dom}(h))$ . Then, we have to show  $\tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}$  from the definition of  $\bar{T}_\gamma^t$ , i.e.,  $\bar{T}_\gamma^t(\mathbf{x})_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}$ . We prove this result by contradiction. Our proof strategy essentially follows the one we developed in Lemma 2 and Lemma 26. For all  $k \in \mathbb{N}$ , we construct a sequence of intermediate points as

$$\mathbf{x}^{\theta,k} := \theta \mathbf{x}^k + (1 - \theta) T_\gamma^t(\mathbf{x}^k),$$

and a sequence of univariate functions as

$$\phi_k(\theta) := \gamma(\mathbf{x}^{\theta,k}; \mathbf{x}) + \frac{1}{t} D_h(\mathbf{x}^{\theta,k}, \mathbf{x}) + g(\mathbf{x}^{\theta,k}),$$

where  $\theta \in (0, 1]$ . Then, we show  $\phi_k(\theta) < \phi_k(0)$  would hold for some  $\theta \in (0, 1]$  and  $k \in \mathbb{N}$  if  $\tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} \neq \mathbf{x}_{\mathcal{B}(\mathbf{x})}$ .

- **Step 2:** We prove  $\tilde{\mathbf{x}} = \lim_{k \rightarrow \infty} T_\gamma^t(\mathbf{x}^k) = \bar{T}_\gamma^t(\mathbf{x})$  via the optimality condition.

For all  $k \in \mathbb{N}$ , from the definition of  $T_\gamma^t(\mathbf{x}^k)$ , we can write down its optimality condition:

$$\nabla \gamma(T_\gamma^t(\mathbf{x}^k); \mathbf{x}^k) + \frac{1}{t} (\nabla h(T_\gamma^t(\mathbf{x}^k)) - \nabla h(\mathbf{x}^k)) + \mathbf{d}^{k+1} = \mathbf{0}, \quad (\text{D.11})$$

where  $\mathbf{d}^{k+1} \in \partial g(T_\gamma^t(\mathbf{x}^k))$ . As  $\mathbf{x}_k \rightarrow \mathbf{x}$  and  $T_\gamma^t(\mathbf{x}^k) \rightarrow \tilde{\mathbf{x}}$ , we can apply Lemma 26 and get the boundedness of  $\{\mathbf{d}^{k+1}\}_{k \in \mathbb{N}}$ . By passing to a subsequence, we can assume that  $\mathbf{d}^{k+1} \rightarrow \bar{\mathbf{d}}$ . Then,

we have  $\bar{\mathbf{d}} \in \partial g(\tilde{\mathbf{x}})$ . As  $k$  approaches infinity in (D.11), and given that the limit of  $\nabla h(\mathbf{x}^k)$  exists for the coordinates corresponding to  $\mathcal{I}(\mathbf{x})$ , it follows that

$$\nabla_{\mathcal{I}(\mathbf{x})}\gamma(\tilde{\mathbf{x}}; \mathbf{x}) + \frac{1}{t} (\nabla_{\mathcal{I}(\mathbf{x})}h(\tilde{\mathbf{x}}) - \nabla_{\mathcal{I}(\mathbf{x})}h(\mathbf{x})) + \bar{\mathbf{d}}_{\mathcal{I}(\mathbf{x})} = \mathbf{0}. \quad (\text{D.12})$$

Moreover, from **Step 1**, we know  $\tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}$  and thus (D.12) is the optimality condition of the problem

$$\min_{\mathbf{y} \in \mathbb{R}^n} G_\gamma^t(\mathbf{y}; \mathbf{x}) = \gamma(\mathbf{y}; \mathbf{x}) + \delta_{\mathbf{y}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}}(\mathbf{y}) + \frac{1}{t} \sum_{i \in \mathcal{I}(\mathbf{x})} D_\varphi(y_i, x_i) + g(\mathbf{y}). \quad (\text{D.13})$$

From Definition 16, we know  $\tilde{\mathbf{x}} = \bar{T}_\gamma^t(\mathbf{x})$ .

**Scenarios (ii):** The key steps essentially follow those developed in **Scenarios (i)**. We highlight the key differences and omit redundant details for simplicity.

To prove **Step 1**, we have to show the sequence  $\{\bar{T}_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$  is bounded, as we did in Lemma 28. We prove this result by contradiction. By passing to a subsequence, we assume that  $\|\bar{T}_\gamma^t(\mathbf{x}^k)\| \rightarrow +\infty$ . Then, we construct a sequence  $\{\mathbf{y}^k\}_{k \in \mathbb{N}} \subseteq \text{int}(\text{dom}(h)) \cap \mathcal{X}$  that satisfying  $\|\mathbf{x}^k - \mathbf{y}^k\| \leq 2^{-k}$  and  $\|T_\gamma^t(\mathbf{y}^k) - \bar{T}_\gamma^t(\mathbf{x}^k)\| \leq 1$ . Such a sequence exists due to the result of **Scenarios (i)** and the existence of interior points. That is, for each  $k \in \mathbb{N}$ , we can always construct a sequence converging to  $\mathbf{x}^k$ , and pick a point  $\mathbf{y}^k$  in the sequence satisfying the conditions. Consequently, we have  $\|T_\gamma^t(\mathbf{y}^k)\| \rightarrow +\infty$  when  $\mathbf{y}^k \rightarrow \mathbf{x}$ . We obtain the contradiction by applying the same arguments in Lemma 28.

From the boundedness of the sequence  $\{\bar{T}_\gamma^t(\mathbf{x}^k)\}_{k \in \mathbb{N}}$ , we can assume  $\bar{T}_\gamma^t(\mathbf{x}^k) \rightarrow \tilde{\mathbf{x}} \in \text{cl}(\text{dom}(h))$  by passing to a subsequence if necessary. For all  $k \in \mathbb{N}$ , we construct a sequence of intermediate points as

$$\bar{\mathbf{x}}^{\theta, k} := \theta \mathbf{x}^k + (1 - \theta) \bar{T}_\gamma^t(\mathbf{x}^k),$$

and a sequence of univariate functions as

$$\bar{\phi}_k(\theta) := G_\gamma^t(\bar{\mathbf{x}}^{\theta, k}; \mathbf{x}^k),$$

where  $\theta \in (0, 1]$ . We can prove that  $\bar{T}_\gamma^t(\mathbf{x})_{\mathcal{B}(\mathbf{x})} = \tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}$  by the same arguments developed in **Step 1** for **Scenarios (i)**.

**Step 2:** Owing to  $\mathbf{x}_k \rightarrow \mathbf{x}$ , WLOG, we can assume that  $\mathcal{I}(\mathbf{x}^k) = \mathcal{I}_0 \subseteq [n]$ . From the definition of  $\bar{T}_\gamma^t(\mathbf{x}^k)$ , we have

$$\nabla_{\mathcal{I}_0}\gamma(\bar{T}_\gamma^t(\mathbf{x}^k); \mathbf{x}^k) + \frac{1}{t} (\nabla_{\mathcal{I}_0}h(\bar{T}_\gamma^t(\mathbf{x}^k)) - \nabla_{\mathcal{I}_0}h(\mathbf{x}^k)) + \mathbf{d}_{\mathcal{I}_0}^{k+1} = \mathbf{0}, \quad (\text{D.14})$$

where  $\mathbf{d}^{k+1} \in \partial g(\bar{T}_\gamma^t(\mathbf{x}^k))$ . Moreover, the sequence  $\{\mathbf{d}_{\mathcal{I}_0}^{k+1}\}_{k \in \mathbb{N}}$  is bounded due to Corollary 27. Then, there always exists a subsequence such that  $\mathbf{d}_{\mathcal{I}_0}^{k+1} \rightarrow \bar{\mathbf{d}}_{\mathcal{I}_0}$  for some  $\bar{\mathbf{d}} \in \mathbb{R}^n$ . As  $k$  approaches infinity in (D.14), and given that the limit of  $\nabla h(\mathbf{x}^k)$  exists for the coordinates corresponding to  $\mathcal{I}(\mathbf{x})$ , it follows that

$$\nabla_{\mathcal{I}(\mathbf{x})}\gamma(\tilde{\mathbf{x}}; \mathbf{x}) + \frac{1}{t} (\nabla_{\mathcal{I}(\mathbf{x})}h(\tilde{\mathbf{x}}) - \nabla_{\mathcal{I}(\mathbf{x})}h(\mathbf{x})) + \bar{\mathbf{d}}_{\mathcal{I}(\mathbf{x})} = \mathbf{0}.$$

Lastly, due to  $\tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} = \mathbf{x}_{\mathcal{B}(\mathbf{x})}$ , we know  $\overline{T}_\gamma^t(\mathbf{x}) = \tilde{\mathbf{x}}$ .

For the final section, our objective is to furnish a comprehensive proof within **Step 1** under **Scenarios (i)**. To demonstrate that  $\phi_k(\theta) < \phi_k(0)$  holds for some  $\theta \in (0, 1]$  and  $k \in \mathbb{N}$  if  $\tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} \neq \mathbf{x}_{\mathcal{B}(\mathbf{x})}$ , we decompose  $\phi_k(\theta) - \phi_k(0)$  following the proof outlined in Lemma 2 and 26. By Mean Value Theorem, we have

$$\begin{aligned} & \frac{\phi_k(\theta) - \phi_k(0)}{\theta} \\ &= \frac{\gamma(\mathbf{x}^{\theta,k}; \mathbf{x}) - \gamma(T_\gamma^t(\mathbf{x}^{\theta,k}); \mathbf{x})}{\theta} + \frac{1}{t}(T_\gamma^t(\mathbf{x}^k) - \mathbf{x}^k)^\top (\nabla h(z^{\theta,k}) - \nabla h(\mathbf{x}^k)) \\ & \quad + \frac{g(\mathbf{x}^{\theta,k}) - g(T_\gamma^t(\mathbf{x}^k))}{\theta}, \end{aligned}$$

where  $z^{\theta,k}$  is between  $\mathbf{x}^{\theta,k}$  and  $T_\gamma^t(\mathbf{x}^k)$ . Following the same argument we did in (D.7). The first term and the third term are uniformly bounded. We just focus on the second term here.

Since we have assumed  $\tilde{\mathbf{x}}_{\mathcal{B}(\mathbf{x})} \neq \mathbf{x}_{\mathcal{B}(\mathbf{x})}$ , there is a  $b_* \in \mathcal{B}(\mathbf{x})$  such that  $\tilde{x}_{b_*} \neq x_{b_*}$ . As  $z_{b_*}^{\theta,k}$  lies in the interval between  $x_{b_*}^{\theta,k}$  and  $T_\gamma^t(\mathbf{x}^k)_{b_*}$ , we have  $z_{b_*}^{\theta,k} - x_{b_*}^{\theta,k} = \xi \cdot (T_\gamma^t(\mathbf{x}^k)_{b_*} - x_{b_*}^k)$  for some  $\xi > 0$ . The monotone of  $\varphi'$  yields

$$\left( \varphi'(z_{b_*}^{\theta,k}) - \varphi'(x_{b_*}^k) \right) \cdot \left( x_{b_*}^k - T_\gamma^t(\mathbf{x}^k)_{b_*} \right) \leq 0. \quad (\text{D.15})$$

Using  $\tilde{x}_{b_*} \neq x_{b_*}$ , we have  $x_{b_*}^k - T_\gamma^t(\mathbf{x}^k)_{b_*} \not\rightarrow 0$ . Then, noticing  $|\varphi'(x_{b_*}^k)| \rightarrow \infty$  due to  $x_{b_*} \in \text{bd}(\text{dom}(\varphi))$ , we have

$$\left( \varphi'(z_{b_*}^{\theta,k}) - \varphi'(x_{b_*}^k) \right) \cdot \left( x_{b_*}^k - T_\gamma^t(\mathbf{x}^k)_{b_*} \right) \rightarrow -\infty \text{ as } k \rightarrow \infty.$$

Thus, we get

$$\lim_{k \rightarrow \infty} \frac{\phi_k(\theta) - \phi_k(0)}{\theta} = -\infty,$$

which yields a contradiction. We complete our proof. ■

## D.2. Proof of Theorem 20

**Proof** For  $\mathbf{x} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , we have  $\mathcal{I}(\mathbf{x}) = [n]$  and the continuity of  $\overline{R}_\gamma^t$  at  $\mathbf{x}$  follows from the continuity of  $\overline{T}_\gamma^t$  at  $\mathbf{x}$ . Then, we only need to show that  $\overline{T}_\gamma^t$  is continuous at  $\mathbf{x} \in \text{bd}(\text{dom}(h))$ . That is, for any  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  converging to  $\mathbf{x} \in \text{bd}(\text{dom}(h))$ ,

$$\lim_{k \rightarrow \infty} \overline{R}_\gamma^t(\mathbf{x}^k) = \overline{R}_\gamma^t(\mathbf{x}). \quad (\text{D.16})$$

Recall the definition of  $\overline{R}_\gamma^t$  as

$$\overline{R}_\gamma^t(\mathbf{x}) = \sum_{i \in \mathcal{I}(\mathbf{x})} D_\varphi(\overline{T}_\gamma^t(\mathbf{x})_i, x_i).$$



Then, the main difficulty to verifying (D.16) is that  $\mathcal{I}(\mathbf{x}^k)$  may not equal  $\mathcal{I}(\mathbf{x})$  for  $k$  sufficiently large.

WLOG, we can assume that  $\mathcal{I}(\mathbf{x}^k) \equiv \mathcal{I}_0 \subseteq [n]$  for all  $k \in \mathbb{N}$ . Then, we have  $\mathcal{I}(\mathbf{x}) \subseteq \mathcal{I}_0$  as  $\mathbf{x}^k \rightarrow \mathbf{x}$ . Now, we discuss  $\mathcal{I}(\mathbf{x})$  and  $\mathcal{I}_0 \setminus \mathcal{I}(\mathbf{x})$  separately.

(i) For all  $i \in \mathcal{I}(\mathbf{x})$ , we have

$$D_\varphi \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i, x_i^k \right) \rightarrow D_\varphi \left( \overline{T}_\gamma^t(\mathbf{x})_i, x_i \right),$$

due to the continuity of  $\overline{T}_\gamma^t$  and  $x_i^k \rightarrow x_i \in \text{int}(\text{dom}(\varphi))$ .

(ii) For all  $i \in \mathcal{I}_0 \setminus \mathcal{I}(\mathbf{x})$ , we want to show

$$D_\varphi \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i, x_i^k \right) \rightarrow 0.$$

To see this, we notice that the convexity of  $\varphi$  yields

$$D_\varphi \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i, x_i^k \right) \leq \left( \varphi' \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i \right) - \varphi'(x_i^k) \right) \cdot \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i - x_i^k \right).$$

To show the right-hand side goes to zero for  $i \in \mathcal{I}_0 \setminus \mathcal{I}(\mathbf{x})$ , we revisit the optimality condition (D.14). Owing to the boundedness of  $\{\mathbf{d}_{\mathcal{I}_0}^{k+1}\}_{k \in \mathbb{N}}$  and  $\nabla \gamma(\overline{T}_\gamma^t(\mathbf{x}^k); \mathbf{x}^k)$ , it follows that

$$\left| \varphi' \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i \right) - \varphi'(x_i^k) \right| \leq \mathcal{O}(1).$$

By continuity of  $\overline{T}_\gamma^t$  and  $\mathbf{x}^k \rightarrow \mathbf{x}$ , we have  $x_i^k \rightarrow x_i$  and  $\overline{T}_\gamma^t(\mathbf{x}^k)_i \rightarrow \overline{T}_\gamma^t(\mathbf{x})_i$  for  $i \in \mathcal{I}_0 \setminus \mathcal{I}(\mathbf{x})$ . Note that  $\mathcal{I}_0 \setminus \mathcal{I}(\mathbf{x}) \subseteq \mathcal{B}(\mathbf{x})$  and  $\overline{T}_\gamma^t(\mathbf{x})_b = x_b$  for  $b \in \mathcal{B}(\mathbf{x})$  by the definition of  $\overline{T}_\gamma^t$ . It follows that  $\overline{T}_\gamma^t(\mathbf{x}^k)_i - x_i^k \rightarrow 0$  for  $i \in \mathcal{I}_0 \setminus \mathcal{I}(\mathbf{x})$ , which completes the proof.

**Remark 29** We remark that the non-symmetry of Bregman divergence, i.e.,  $D_h(\mathbf{y}, \mathbf{x}) \neq D_h(\mathbf{x}, \mathbf{y})$ , does not affect our results. For instance, if  $R^t(\mathbf{x}) = D_h(\mathbf{x}, T_\gamma^t(\mathbf{x}))$ , and correspondingly,

$$\overline{R}_\gamma^t(\mathbf{x}) := \sum_{i \in \mathcal{I}(\mathbf{x})} D_\varphi \left( x_i, \overline{T}_\gamma^t(\mathbf{x})_i \right),$$

then we still get the continuity of  $\overline{R}_\gamma^t$  by applying the proof of Theorem 20 with Bregman divergences  $D_\varphi \left( \overline{T}_\gamma^t(\mathbf{x}^k)_i, x_i^k \right)$  and  $D_\varphi \left( \overline{T}_\gamma^t(\mathbf{x})_i, x_i \right)$  replaced by  $D_\varphi \left( x_i^k, \overline{T}_\gamma^t(\mathbf{x}^k)_i \right)$  and  $D_\varphi \left( x_i, \overline{T}_\gamma^t(\mathbf{x})_i \right)$ . ■

## Appendix E. Proof of main results

In this section, we provide the missing proofs of main results in Sec. 3, which are mainly based on the continuity property of the extended stationarity measure and the definition of spurious stationary points. We note that the proof of Proposition 5 has been given in Sec. C.4.

### E.1. Proof of Theorem 7

**Proof** For any  $K \in \mathbb{N}$  and arbitrary  $\epsilon > 0$ , our goal is to construct the initial point  $\mathbf{x}^0$  sufficiently close to the spurious point  $\tilde{\mathbf{x}}$  such that  $\|\mathbf{x}^K - \tilde{\mathbf{x}}^*\| \leq \epsilon$ .

The key step lies in proving the following claim is correct: Given an arbitrary  $\epsilon_0 > 0$ , there exist  $\epsilon_1 \in (0, \frac{1}{2}\epsilon_0]$  such that

$$\|\mathbf{x}^K - \tilde{\mathbf{x}}^*\| \leq \epsilon_0,$$

whenever  $\|\mathbf{x}^{K-1} - \tilde{\mathbf{x}}^*\| \leq \epsilon_1$ . At first, we have  $T_\gamma^t(\mathbf{x}^{K-1}) = \bar{T}_\gamma^t(\mathbf{x}^{K-1})$  since  $\mathbf{x}^{K-1}$  is the interior point and Proposition 25, we have

$$\begin{aligned} \|\mathbf{x}^K - \tilde{\mathbf{x}}^*\| &= \|T_\gamma^t(\mathbf{x}^{K-1}) - \bar{T}_\gamma^t(\mathbf{x}^{K-1}) + \bar{T}_\gamma^t(\mathbf{x}^{K-1}) - \bar{T}_\gamma^t(\tilde{\mathbf{x}}^*) + \bar{T}_\gamma^t(\tilde{\mathbf{x}}^*) - \tilde{\mathbf{x}}^*\|_2 \\ &= \|\bar{T}_\gamma^t(\mathbf{x}^{K-1}) - \bar{T}_\gamma^t(\tilde{\mathbf{x}}^*)\|_2. \end{aligned}$$

Moreover, due to Theorem 20, we know the mapping  $\bar{T}_\gamma^t$  is continuous. Thus, for any  $\epsilon_0 > 0$ , there exists some constants  $\delta > 0$  such that  $\|\mathbf{x}^{K-1} - \tilde{\mathbf{x}}^*\|_2 < \delta$ , we have  $\|\bar{T}_\gamma^t(\mathbf{x}^{K-1}) - \bar{T}_\gamma^t(\tilde{\mathbf{x}}^*)\|_2 < \epsilon_0$ . We can always choose a small constant  $\epsilon' < \min\{\delta, \frac{1}{2}\epsilon_0\}$  to make the above argument hold.

Repeating the above argument  $K$  times, there exists a sequence  $\{\epsilon_k\}_{k=0}^K$  such that  $\epsilon_{k+1} \leq \frac{1}{2}\epsilon_k$  and

$$\|\mathbf{x}^{K-k} - \tilde{\mathbf{x}}^*\| \leq \epsilon_k,$$

whenever we have  $\|\mathbf{x}^{K-k-1} - \tilde{\mathbf{x}}^*\| \leq \epsilon_{k+1}$ , for any  $k \in [K]$ .

Now, we are ready to construct the initial point, i.e.,  $\mathbf{x}^0 \in \mathbb{B}_{\epsilon_K}(\tilde{\mathbf{x}}^*) \cap \mathcal{X} \cap \text{int}(\text{dom}(h))$ . We set  $\epsilon_0 = \epsilon$  and get

$$\mathbf{x}^k \in \mathbb{B}_{\epsilon_{K-k}}(\tilde{\mathbf{x}}^*) \subseteq \mathbb{B}_\epsilon(\tilde{\mathbf{x}}^*) \text{ for } k = 0, 1, \dots, K.$$

We complete our proof. ■

### E.2. Proof of Proposition 9

**Proof** Since  $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^n$ , then  $\text{cl}(\text{dom}(\varphi)) = \mathbb{R}_+$ ,  $\mathcal{I}(\mathbf{x}) = \{i \in [n] : x_i > 0\}$ , and  $\mathcal{B}(\mathbf{x}) = \{b \in [n] : x_b = 0\}$ . Moreover, we have  $g(\mathbf{x}) = \delta_{\mathcal{X}}(\mathbf{x})$  and

$$\partial g(\mathbf{x}) = \{A^T \boldsymbol{\mu} - \boldsymbol{\lambda} : \boldsymbol{\mu} \in \mathbb{R}^m, \lambda_i = 0 \forall i \in \mathcal{I}(\mathbf{x}), \lambda_b \geq 0 \forall b \in \mathcal{B}(\mathbf{x})\}.$$

The compactness of  $\mathcal{X}$  ensures the existence of  $\tilde{\mathbf{x}}^* = \text{argmax}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . Since  $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$  is convex problem and  $f$  is not a constant on  $\mathcal{X}$ , we have  $f(\tilde{\mathbf{x}}^*) \neq \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . Thus, from optimality condition, we have

$$0 \notin \nabla f(\tilde{\mathbf{x}}^*) + \partial g(\tilde{\mathbf{x}}^*). \quad (\text{E.1})$$

By contrast,  $\tilde{\mathbf{x}}^*$  is the optimal solution of problem  $\min_{\mathbf{x} \in \mathcal{X}} -f(\mathbf{x})$ , whose optimality condition yields  $\mathbf{0} \in -\nabla f(\tilde{\mathbf{x}}^*) + \partial \delta_{\mathcal{X}}(\tilde{\mathbf{x}}^*)$ , which is equivalent to

$$\mathbf{0} \in \{-\nabla f(\tilde{\mathbf{x}}^*) + A^T \boldsymbol{\mu} - \boldsymbol{\lambda} : \boldsymbol{\mu} \in \mathbb{R}^m, \lambda_i = 0 \forall i \in \mathcal{I}(\tilde{\mathbf{x}}^*), \lambda_b \geq 0 \forall b \in \mathcal{B}(\tilde{\mathbf{x}}^*)\}$$

It follows that there exists  $\boldsymbol{\mu}$  and  $\boldsymbol{\lambda}$  such that

$$\boldsymbol{\lambda}_{\mathcal{B}(\tilde{\mathbf{x}}^*)} \geq \mathbf{0}, \boldsymbol{\lambda}_{\mathcal{I}(\tilde{\mathbf{x}}^*)} = \mathbf{0}, \text{ and } \mathbf{0} = -\nabla f(\tilde{\mathbf{x}}^*) + A^T \boldsymbol{\mu} - \boldsymbol{\lambda}.$$

Let  $\mathbf{p} = -\boldsymbol{\lambda}$ . Then, we have

$$\mathbf{p} = \nabla f(\tilde{\mathbf{x}}^*) - A^T \boldsymbol{\mu} \in \nabla f(\tilde{\mathbf{x}}^*) + \partial g(\tilde{\mathbf{x}}^*) \text{ and } \mathbf{p}_{\mathcal{I}(\tilde{\mathbf{x}}^*)} = 0.$$

Together with (E.1), we can conclude that  $\tilde{\mathbf{x}}^*$  is a spurious point. ■

## Appendix F. Examples of spurious stationary points

We present simple counter-examples, both convex and non-convex, to illustrate the presence of spurious stationary points.

**Example 3 (Convex counter-example)** *Suppose that  $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^2$  and consider the following simple problem:*

$$\begin{aligned} \min_{x_1, x_2} \quad & -x_1 \\ \text{s.t.} \quad & x_1 + x_2 = 1, x_1, x_2 \geq 0. \end{aligned}$$

The point  $(0, 1)$  is identified as a spurious stationary point. We can determine the interior coordinate  $\mathcal{I}((0, 1)) = 2$  and compute the subdifferential at the point  $(0, 1)$  as

$$\begin{aligned} \partial F((0, 1)) &= (-1, 0) + \mathcal{N}_{\{\mathbf{x} \in \mathbb{R}_+^2 : x_1 + x_2 = 1\}}((0, 1)) \\ &= \{(-1, 0) + \lambda(-1, 0) + \mu(1, 1) : \lambda \in \mathbb{R}_+, \mu \in \mathbb{R}\}. \end{aligned}$$

Consequently, we find that  $\mathbf{0} \notin \partial F((0, 1))$  and  $\mathbf{p} = (-1, 0) \in \partial F((0, 1))$  with  $\mathbf{p}_{\mathcal{I}((0, 1))} = p_2 = 0$ .

**Remark 30** *It is worth noting that in Example 3, for all feasible points  $\mathbf{x}$  lying in the interior,*

$$\text{dist}(\mathbf{0}, \partial F(\mathbf{x})) = \min_{\mu \in \mathbb{R}} \|(-1, 0) + \mu(1, 1)\| = \frac{\sqrt{2}}{2}.$$

Since a sequence  $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$  generated by the BPs belong to the interior of the kernel domain, we see that  $\text{dist}(\mathbf{0}, \partial F(\mathbf{x}^k)) \equiv 0$  for all  $k \in \mathbb{N}$ . Hence, the minimal subdifferential norm  $\text{dist}(\mathbf{0}, \partial F)$  is not a suitable measure for the BPs.

**Example 4 (Nonconvex counter-example)** *Suppose that  $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^2$  and consider the following simple problem:*

$$\begin{aligned} \min_{x_1, x_2} \quad & f(x_1, x_2) = -x_1^2 + x_2 \\ \text{s.t.} \quad & x_1 + x_2 = 1, x_1, x_2 \geq 0. \end{aligned}$$

Similar with the convex case, the point  $(0, 1)$  is identified as a spurious stationary point. We can determine that the interior coordinate  $\mathcal{I}((0, 1)) = 2$  and compute the subdifferential at the point  $(0, 1)$  as

$$\partial F((0, 1)) = \{(0, 1) + \lambda(-1, 0) + \mu(1, 1) : \lambda \in \mathbb{R}_+, \mu \in \mathbb{R}\}.$$

Consequently, we find that  $\mathbf{0} \notin \partial F((0, 1))$  and  $\mathbf{p} = (-1, 0) \in \partial F((0, 1))$  with  $\mathbf{p}_{\mathcal{I}((0, 1))} = p_2 = 0$ .

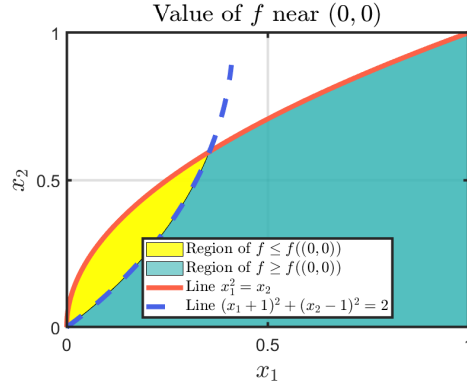


Figure 1: Illustration for Example 5

One may wonder whether spurious stationary points must be (locally) maximal points. The following example reveals that this is not true.

**Example 5** Consider the convex problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}_+^2} \quad & f(x_1, x_2) = (x_1 + 1)^2 + (x_2 - 1)^2 \\ \text{s.t.} \quad & x_2^2 \leq x_1. \end{aligned}$$

The point  $(0, 0)$  is identified as a spurious stationary point. We can determine the interior coordinate  $\mathcal{I}((0, 0)) = \emptyset$  and compute the subdifferential at the point  $(0, 0)$  as

$$\begin{aligned} \partial F((0, 0)) &= (2, -2) + \mathcal{N}_{\{\mathbf{x} \in \mathbb{R}_+^2 : x_2^2 \leq x_1\}}((0, 0)) \\ &= \{(2, -2) + \lambda_1(-1, 0) + \lambda_2(-1, 0) + \lambda_3(0, -1) : \boldsymbol{\lambda} \in \mathbb{R}_+^3\}. \end{aligned}$$

Consequently, we find that  $\mathbf{0} \notin \partial F((0, 0))$ , and hence  $(0, 0)$  is a spurious stationary point.

In Example 5, the spurious stationary point  $(0, 0)$  is not a (local) maximal point. We illustrate this fact through Figure 1. Clearly, the region where  $F \leq F((0, 0))$ , i.e., the yellow part that satisfies  $(x_1 + 1)^2 + (x_2 - 1)^2 \leq 2$ ,  $x_2^2 \leq x_1$ , and  $\mathbf{x} \in \mathbb{R}_+^2$ , is non-empty.

Finally, to enhance our understanding of the finite step trap behavior nearby spurious points, we give several instances on escaping from spurious stationary points.

**Example 6** We revisited the convex counter-example presented in Example 3, where the unique spurious point is  $\tilde{\mathbf{x}}^* = (0, 1)$ . We choose the kernel function as the negative entropy  $\varphi(x) = x \log(x)$ , a popular choice for managing simplex constraints. For any  $K \in \mathbb{N}$  and  $\epsilon > 0$ , we construct the initial point  $\mathbf{x}^0$  as follows:

$$\mathbf{x}^0 = \left( \frac{\sqrt{2}\epsilon}{2} e^{-tK}, 1 - \frac{\sqrt{2}\epsilon}{2} e^{-tK} \right).$$

Moreover, for all  $k \in [K]$ , we apply the standard Bregman gradient descent method as

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{y}} t(-1, 0)^T \mathbf{y} + D_h(\mathbf{y}, \mathbf{x}^k) + \delta_{\Delta_2}(\mathbf{y})$$

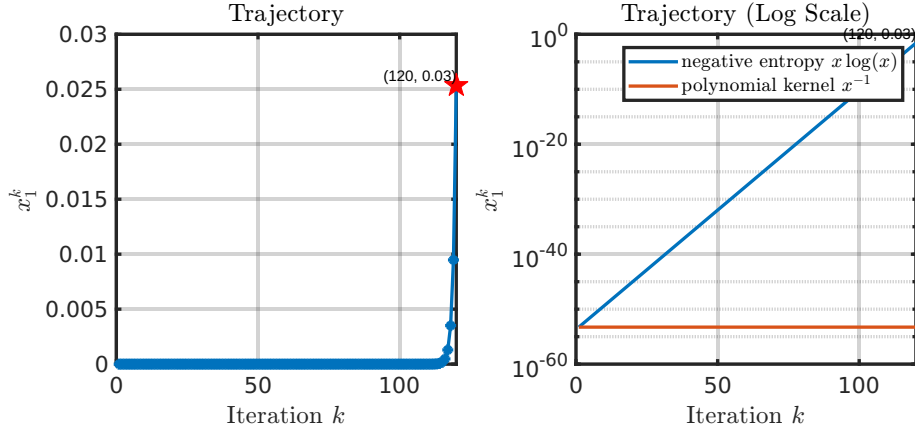


Figure 2: The trajectory plot for  $\{x_1^k\}_{k \in [K]}$  in Example 6 and 7 ( $\alpha = 1$ ), with  $K = 120$  and  $\epsilon = 0.1$ . The initial point is chosen according to the negative entropy kernel scenario.

$$= \left( \frac{x_1^k}{x_1^k + e^{-t}x_2^k}, \frac{e^{-t}x_2^k}{x_1^k + e^{-t}x_2^k} \right), \quad \forall k \in [K].$$

Then, we can quantify the distance between  $\mathbf{x}^{k+1}$  and  $\tilde{\mathbf{x}}^*$  as

$$\|\mathbf{x}^{k+1} - \tilde{\mathbf{x}}^*\| = \frac{\sqrt{2}x_1^k}{x_1^k + e^{-t}x_2^k} \leq \sqrt{2}e^t x_1^k \leq \sqrt{2}e^{tk} x_1^0 = e^{-t(K-k)} \epsilon \leq \epsilon,$$

where the first inequality is derived from the constraint  $x_1^k + x_2^k = 1$  and  $t \geq 0$ , the second inequality is justified by iteratively applying the recursive relation from the first inequality  $k$  times.

From the example involving negative entropy, it becomes clear that constructing the initial point at a distance that exponentially decays with respect to the spurious point is crucial. Later on, we want to provide another artificially constructed example to demonstrate the importance of the kernel function's growth condition in determining the necessary distance between the initial point and the spurious point to trigger a finite step trap. Essentially, the challenge of falling into a finite step trap varies significantly across different kernel functions.

**Example 7 (Polynomial kernel)** We still consider the convex counter-example presented in Example 3 with a kernel function as  $\varphi(x) = \frac{1}{\alpha}x^{-\alpha}$  where  $\alpha > 0$ . For any  $K \in \mathbb{N}$  and  $\epsilon > 0$ , we construct the initial point  $\mathbf{x}^0$  as follows:

$$\mathbf{x}_1^0 = \min \left\{ \left( \frac{2}{tK + \epsilon^{-(\alpha+1)}} \right)^{\frac{1}{\alpha+1}}, \frac{1}{1 + 2^{\alpha+1}} \right\}$$

and  $x_2^0 = 1 - x_1^0$ . Moreover, for all  $k \in [K]$ , we apply the standard Bregman gradient descent method as

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{y}} t(-1, 0)^T \mathbf{y} + D_h(\mathbf{y}, \mathbf{x}^k) + \delta_{\Delta_2}(\mathbf{y}).$$

If  $\mathbf{x}^k \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$ , we have  $\mathbf{x}^{k+1} \in \text{int}(\text{dom}(h)) \cap \mathcal{X}$  from Lemma 2. We can write down its optimality condition:

$$\begin{cases} -t + \varphi'(x_1^{k+1}) - \varphi'(x_1^k) + \mu_{k+1} = 0 \\ \varphi'(x_2^{k+1}) - \varphi'(x_2^k) + \mu_{k+1} = 0 \\ x_1^{k+1} + x_2^{k+1} = 1. \end{cases}$$

Summing up the above equation from  $k = 0$  to  $k = K - 1$ , we have

$$\begin{cases} -tK + \varphi'(x_1^K) - \varphi'(x_1^0) + \sum_{k=1}^K \mu_{k+1} = 0 \\ \varphi'(x_2^K) - \varphi'(x_2^0) + \sum_{k=1}^K \mu_{k+1} = 0. \end{cases}$$

Since  $\varphi'(x) = -x^{-\alpha-1}$  and  $(\varphi')^{-1}(y) = \left(-\frac{1}{y}\right)^{\frac{1}{1+\alpha}}$ , we know  $\varphi'$  is negative on  $\mathbb{R}_{++}$  and  $(\varphi')^{-1}$  is monotonically increasing on  $\mathbb{R}_-$ . Then, we have

$$-\sum_{k=1}^K \mu_{k+1} = \varphi'(x_2^K) - \varphi'(x_2^0) \leq -\varphi'(x_2^0).$$

Finally, we proceed to bound  $x_1^K$

$$\begin{aligned} x_1^K &= (\varphi')^{-1} \left( \varphi'(x_1^0) - \sum_{k=1}^K \mu_{k+1} + tK \right) \\ &\leq (\varphi')^{-1} \left( \varphi'(x_1^0) - \varphi'(1 - x_1^0) + tK \right) \\ &= \left( (x_1^0)^{-\alpha-1} - (1 - x_1^0)^{-\alpha-1} - tK \right)^{\frac{-1}{1+\alpha}} \leq \epsilon. \end{aligned}$$

where the first inequality follows from  $-\sum_{k=1}^K \mu_{k+1} \leq -\varphi'(x_2^0)$  and  $x_1^0 + x_2^0 = 1$  and the last inequality follows from

$$(x_1^0)^{-\alpha-1} - (1 - x_1^0)^{-\alpha-1} \geq (2x_1^0)^{-\alpha-1}$$

due to  $x_1^0 \leq \frac{1}{1+2^{\alpha+1}}$  and  $(2x_1^0)^{-\alpha-1} - tK > 0$ . For simplicity, we ignore the final constant  $\sqrt{2}$ .

**Remark 31** From Figure 2, it is evident that despite both kernels facing the challenges outlined in Theorem 7, the nonnegative entropy kernel demonstrates superior performance compared to the polynomial kernel. This advantage can be attributed to the curvature information encapsulated by the inverse mapping of  $\nabla h(\mathbf{x})$ . Specifically, when  $h(\mathbf{x}) = \mathbf{x} \log(\mathbf{x})$ ,  $\nabla^{-1}h(\mathbf{x})$  exhibits exponential growth behavior, enabling each iteration to escape the unfavorable point by at least doubling the distance from it.