

Spurious Stationarity and Hardness Results for Bregman Proximal-type Algorithms

Jiajin Li

Sauder School of Business

University of British Columbia



September 21, 2024

Joint work with He Chen (CUHK) and Anthony Man-Cho So (CUHK).

1. Spurious stationary points inevitably exist when **non-gradient Lipschitz kernels** are used for Bregman proximal-type algorithms.

2. (**Algorithm-Dependent Hardness Results**) Bregman proximal-type algorithms *are unable to* escape from a spurious stationary point in finite steps when the initial point is bad.

1. Spurious stationary points inevitably exist when **non-gradient Lipschitz kernels** are used for Bregman proximal-type algorithms.
2. (**Algorithm-Dependent Hardness Results**) Bregman proximal-type algorithms *are unable to* escape from a spurious stationary point in finite steps when the initial point is bad.

What we'll cover today?

1. **Introduction and Problem Settings**
2. Spurious Stationary Points and Examples
3. Algorithm-Dependent Hardness Results and their Implications
4. Unsatisfactory Stationary Measures and Convergence Behaviour Investigation

Mirror Descent (Non-Euclidean Gradient Descent)

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \quad (\text{P})$$

- Gradient Descent:

$$\begin{aligned} \mathbf{x}_+ &= \mathbf{x} - t \cdot \nabla F(\mathbf{x}) \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^n} \nabla F(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

- Mirror Descent:

$$\mathbf{x}_+ = \arg \min_{\mathbf{y} \in \mathbb{R}^n} \nabla F(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2t} \underbrace{D_h(\mathbf{y}, \mathbf{x})}_{\text{Bregman Divergence}}.$$

Better to exploit the geometry of the problem at hand!

Mirror Descent (Non-Euclidean Gradient Descent)

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \quad (\text{P})$$

- Gradient Descent:

$$\begin{aligned} \mathbf{x}_+ &= \mathbf{x} - t \cdot \nabla F(\mathbf{x}) \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^n} \nabla F(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

- Mirror Descent:

$$\mathbf{x}_+ = \arg \min_{\mathbf{y} \in \mathbb{R}^n} \nabla F(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2t} \underbrace{D_h(\mathbf{y}, \mathbf{x})}_{\text{Bregman Divergence}}.$$

Better to exploit the geometry of the problem at hand!

Mirror Descent (Non-Euclidean Gradient Descent)

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \quad (\text{P})$$

- Gradient Descent:

$$\begin{aligned} \mathbf{x}_+ &= \mathbf{x} - t \cdot \nabla F(\mathbf{x}) \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^n} \nabla F(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

- Mirror Descent:

$$\mathbf{x}_+ = \arg \min_{\mathbf{y} \in \mathbb{R}^n} \nabla F(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2t} \underbrace{D_h(\mathbf{y}, \mathbf{x})}_{\text{Bregman Divergence}}.$$

Better to exploit the geometry of the problem at hand!

Definition

The Bregman divergence between two points \mathbf{x}, \mathbf{y} associated with a kernel function $h : \Omega \rightarrow \mathbb{R}$ is defined as

$$D_h(\mathbf{x}, \mathbf{y}) := h(\mathbf{x}) - h(\mathbf{y}) - \nabla h(\mathbf{y})^T (\mathbf{x} - \mathbf{y}),$$

where h is continuously differentiable and strictly convex on the convex set Ω .

- (✓) If $h(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$, we have $D_h(\mathbf{y}, \mathbf{x}) = \frac{1}{2}\|\mathbf{y} - \mathbf{x}\|^2$. Then, mirror descent reduces to the vanilla gradient descent.
- (✗) If $h(\mathbf{x}) = \sum_{i=1}^n x_i \log x_i$, $D_h(\mathbf{y}, \mathbf{x})$ is just KL divergence.
- (✗) Other non-gradient Lipschitz kernel functions ...

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

- $\text{dom}(g) = \mathcal{X}$ is a nonempty closed convex set.
- $f : \mathcal{X} \rightarrow \mathbb{R}$ is continuous differentiable on \mathcal{X} (possibly **nonconvex**).
- $g : \mathcal{X} \rightarrow \overline{\mathbb{R}}$ is **convex** and **locally Lipschitz continuous**, e.g.,
 - Indicator function \rightarrow include the constrained optimization problem as a special case.

Bregman Proximal-Type Algorithms

$$\mathbf{x}_+ = T_\gamma^t(\mathbf{x}) := \arg \min_{\mathbf{y} \in \mathbb{R}^n} \left\{ \underbrace{\gamma(\mathbf{y}; \mathbf{x})}_{\text{Surrogate Model}} + g(\mathbf{y}) + \frac{1}{t} D_h(\mathbf{y}, \mathbf{x}) \right\} \quad (\mathcal{A})$$

- If $\gamma(\mathbf{y}; \mathbf{x}) = f(\mathbf{y})$, (\mathcal{A}) reduces to Bregman proximal point methods [Chen and Teboulle, 1993].
- If $\gamma(\mathbf{y}; \mathbf{x}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$, (\mathcal{A}) reduces to Bregman proximal (projected) gradient descent [Bauschke et al., 2017],[Bauschke et al., 2019].
- If $\gamma(\mathbf{y}; \mathbf{x}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})$, (\mathcal{A}) has been recently explored by [Doikov and Nesterov, 2023].

Separable Kernel Functions

- $h(\mathbf{x}) = \sum_{i=1}^n \varphi(x_i)$, where $\varphi : \mathbb{R} \rightarrow \overline{\mathbb{R}}$ is a univariate function.
- φ is continuously differentiable on $\text{int}(\text{dom}(\varphi))$, and $|\varphi'(x^k)| \rightarrow +\infty$ as $x^k \rightarrow x \in \text{bd}(\text{dom}(\varphi))$.
- φ is strictly convex.

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

Separable Kernel Functions

- $h(\mathbf{x}) = \sum_{i=1}^n \varphi(x_i)$, where $\varphi : \mathbb{R} \rightarrow \overline{\mathbb{R}}$ is a univariate function.
- φ is continuously differentiable on $\text{int}(\text{dom}(\varphi))$, and $|\varphi'(x^k)| \rightarrow +\infty$ as $x^k \rightarrow x \in \text{bd}(\text{dom}(\varphi))$.
- φ is strictly convex.

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

What we'll cover today?

1. Introduction and Problem Settings
2. **Spurious Stationary Points and Examples**
3. Algorithm-Dependent Hardness Results and their Implications
4. Unsatisfactory Stationary Measures and Convergence Behaviour Investigation

Definition

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 $0 \notin \partial F(\mathbf{x})$.

- $\mathcal{I}(\mathbf{x}) := \{i \in [n] : x_i \in \text{int}(\text{dom}(\varphi))\}$.
- Spurious stationary points exist only when the kernel is **non-gradient Lipschitz**.
- For a kernel h with gradient Lipschitz property, we have $\text{dom}(\varphi) = \mathbb{R}$ and $\mathcal{I}(\mathbf{x}) = [n]$ hold for all $\mathbf{x} \in \mathcal{X}$, thereby precluding the existence of spurious stationary points.

Only depends on the problem itself and the kernel function!

Definition

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 $0 \notin \partial F(\mathbf{x})$.

- $\mathcal{I}(\mathbf{x}) := \{i \in [n] : x_i \in \text{int}(\text{dom}(\varphi))\}$.
- Spurious stationary points exist only when the kernel is **non-gradient Lipschitz**.
- For a kernel h with gradient Lipschitz property, we have $\text{dom}(\varphi) = \mathbb{R}$ and $\mathcal{I}(\mathbf{x}) = [n]$ hold for all $\mathbf{x} \in \mathcal{X}$, thereby precluding the existence of spurious stationary points.

Only depends on the problem itself and the kernel function!

Spurious Stationarity

Definition

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 $0 \notin \partial F(\mathbf{x})$.

- $\mathcal{I}(\mathbf{x}) := \{i \in [n] : x_i \in \text{int}(\text{dom}(\varphi))\}$.
- Spurious stationary points exist only when the kernel is **non-gradient Lipschitz**.
- For a kernel h with gradient Lipschitz property, we have $\text{dom}(\varphi) = \mathbb{R}$ and $\mathcal{I}(\mathbf{x}) = [n]$ hold for all $\mathbf{x} \in \mathcal{X}$, thereby precluding the existence of spurious stationary points.

Only depends on the problem itself and the kernel function!

Example (A Simple Linear Programming Problem)

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 find that $0 \notin \partial F((0,1))$ and $p = (-1, 0) \in \partial F((0,1))$ with $p_{\mathcal{I}((0,1))} = p_2 = 0$, i.e.,
$$\partial F((0,1)) = \{(-1, 0) + \lambda(-1, 0) + \mu(1, 1) : \lambda \in \mathbb{R}_+, \mu \in \mathbb{R}\}.$$

Example (A Simple Linear Programming Problem)

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

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

Nonconvex Example

Example

Suppose that $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^2$ and consider the following simple problem:

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

The point $(0,1)$ is identified as a spurious stationary point.

Proposition (Existence of Spurious Stationary Points)

Consider a convex optimization problem

$$\begin{aligned} \min_{x \in \mathcal{X}} \quad & f(x) \\ \text{s.t.} \quad & Ax = \mathbf{b}, x \geq 0. \end{aligned}$$

Suppose the constraint set is compact and f is non-constant. If $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^n$, then every maximal point $\tilde{x}^ \in \text{argmax}_{x \in \mathcal{X}} f(x)$ is a spurious stationary point.*

Nonconvex Example

Example

Suppose that $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^2$ and consider the following simple problem:

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

The point $(0,1)$ is identified as a spurious stationary point.

Proposition (Existence of Spurious Stationary Points)

Consider a convex optimization problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0. \end{aligned}$$

*Suppose the constraint set is compact and f is non-constant. If $\text{cl}(\text{dom}(h)) = \mathbb{R}_+^n$, then **every maximal point** $\tilde{\mathbf{x}}^* \in \text{argmax}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ is a spurious stationary point.*

What we'll cover today?

1. Introduction and Problem Settings
2. Spurious Stationary Points and Examples
3. **Algorithm-Dependent Hardness Results and their Implications**
4. Unsatisfactory Stationary Measures and Convergence Behaviour Investigation

Algorithm-dependent Hardness Result

Theorem

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 $\mathbf{x}^0 \in \mathcal{B}_\epsilon(\tilde{\mathbf{x}}^*) \cap \mathcal{X}$, sufficiently close to the **spurious stationary point** $\tilde{\mathbf{x}}^*$, such that

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

[Corollary 1, Bauschke et al., 2017]. When f is convex, the sequence $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ generated by BPG satisfies

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

where $\bar{\mathbf{x}} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f$ is the global minimizer, t is the step size, and \mathbf{x}^0 is an arbitrary initial point.

Algorithm-dependent Hardness Result

Theorem

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 $\mathbf{x}^0 \in \mathcal{B}_\epsilon(\tilde{\mathbf{x}}^*) \cap \mathcal{X}$, sufficiently close to the **spurious stationary point** $\tilde{\mathbf{x}}^*$, such that

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

[Corollary 1, Bauschke et al., 2017]. When f is convex, the sequence $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ generated by BPG satisfies

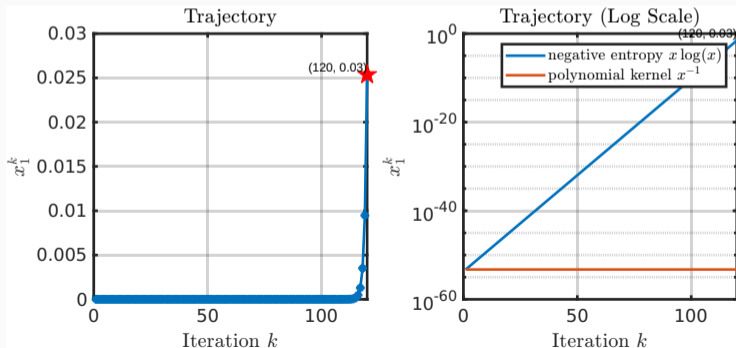
$$f(\mathbf{x}^k) - \min_{\mathbf{x} \in \mathcal{X}} f \leq \overbrace{\frac{D_h(\bar{\mathbf{x}}, \mathbf{x}^0)}{t}}^{\text{Extremely Large}} \cdot \frac{1}{k},$$

where $\bar{\mathbf{x}} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f$ is the global minimizer, t is the step size, and \mathbf{x}^0 is an arbitrary initial point.

Example (The Simple Linear Programming Problem)

For every K and $\epsilon > 0$, we construct the initial point as

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



What we'll cover today?

1. Introduction and Problem Settings
2. Spurious Stationary Points and Examples
3. Algorithm-Dependent Hardness Results and their Implications
4. **Unsatisfactory Stationary Measures and Convergence Behaviour**

Understand how the sequence of iterations behaves and how close it gets to convergence.

A standard recipe in optimization:

- Propose a residual function $R : \mathbb{R}^n \rightarrow \mathbb{R}_+$ that measures the stationarity of the iterations.
- Establish the convergence of the sequence of $\{R(\mathbf{x}^k)\}_{k \geq 0}$.

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

Understand how the sequence of iterations behaves and how close it gets to convergence.

A standard recipe in optimization:

- Propose a residual function $R : \mathbb{R}^n \rightarrow \mathbb{R}_+$ that measures the stationarity of the iterations.
- Establish the convergence of the sequence of $\{R(\mathbf{x}^k)\}_{k \geq 0}$.

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

Understand how the sequence of iterations behaves and how close it gets to convergence.

A standard recipe in optimization:

- Propose a residual function $R : \mathbb{R}^n \rightarrow \mathbb{R}_+$ that measures the stationarity of the iterations.
- Establish the convergence of the sequence of $\{R(\mathbf{x}^k)\}_{k \geq 0}$.

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

Existing stationarity measures are not well-defined

All existing stationarity measure can be unified as

$$R_\gamma^t(\mathbf{x}) := D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})$$

the relative change w.r.t Bregman divergence.

- R_γ^t is not well-defined on the boundary $\text{bd}(\text{dom}(h))$.
- The mapping $\mathbf{x} \mapsto T_\gamma^t(\mathbf{x})$ involves the Bregman divergence function $(\mathbf{y}, \mathbf{x}) \mapsto D_h(\mathbf{y}, \mathbf{x})$, which is only defined on $\text{dom}(h) \times \text{int}(\text{dom}(h))$.

Existing stationarity measures are not well-defined

All existing stationarity measure can be unified as

$$R_\gamma^t(\mathbf{x}) := D_h(T_\gamma^t(\mathbf{x}), \mathbf{x})$$

the relative change w.r.t Bregman divergence.

- R_γ^t is not well-defined on the boundary $\text{bd}(\text{dom}(h))$.
- The mapping $\mathbf{x} \mapsto T_\gamma^t(\mathbf{x})$ involves the Bregman divergence function $(\mathbf{y}, \mathbf{x}) \mapsto D_h(\mathbf{y}, \mathbf{x})$, which is only defined on $\text{dom}(h) \times \text{int}(\text{dom}(h))$.

Extended Stationarity Measure

A simple fix: Only account for the **interior coordinates**, i.e.,

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

where $\bar{T}_\gamma^t(\mathbf{x})$ denotes the update rule that ensures the boundary coordinates remain fixed.

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

- The residual function is **continuous** (✓).
- The residual function equals to zeros **if and only** if \mathbf{x} is a stationary point.
 - If \mathbf{x} is a stationary point, we have $\bar{R}_\gamma^t(\mathbf{x}) = 0$ (✓).
 - $\Rightarrow ?$

Extended Stationarity Measure

A simple fix: Only account for the **interior coordinates**, i.e.,

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

where $\bar{T}_\gamma^t(\mathbf{x})$ denotes the update rule that ensures the boundary coordinates remain fixed.

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

- The residual function is **continuous** (✓).
- The residual function equals to zeros **if and only** if \mathbf{x} is a stationary point.
 - If \mathbf{x} is a stationary point, we have $\bar{R}_\gamma^t(\mathbf{x}) = 0$ (✓).
 - $\Rightarrow ?$

Extended Stationarity Measure

A simple fix: Only account for the **interior coordinates**, i.e.,

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

where $\bar{T}_\gamma^t(\mathbf{x})$ denotes the update rule that ensures the boundary coordinates remain fixed.

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

- The residual function is **continuous** (✓).
- The residual function equals to zeros **if and only** if \mathbf{x} is a stationary point.
 - If \mathbf{x} is a stationary point, we have $\bar{R}_\gamma^t(\mathbf{x}) = 0$ (✓).
 - $\Rightarrow ?$

Extended Stationarity Measure

A simple fix: Only account for the **interior coordinates**, i.e.,

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

where $\bar{T}_\gamma^t(\mathbf{x})$ denotes the update rule that ensures the boundary coordinates remain fixed.

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

- The residual function is **continuous** (✓).
- The residual function equals to zeros **if and only** if \mathbf{x} is a stationary point.
 - If \mathbf{x} is a stationary point, we have $\bar{R}_\gamma^t(\mathbf{x}) = 0$ (✓).

• $\Rightarrow ?$

Extended Stationarity Measure

A simple fix: Only account for the **interior coordinates**, i.e.,

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

where $\bar{T}_\gamma^t(\mathbf{x})$ denotes the update rule that ensures the boundary coordinates remain fixed.

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

- The residual function is **continuous** (✓).
- The residual function equals to zeros **if and only** if \mathbf{x} is a stationary point.
 - If \mathbf{x} is a stationary point, we have $\bar{R}_\gamma^t(\mathbf{x}) = 0$ (✓).
 - $\Rightarrow ?$

Proposition (Characterization of spurious stationary points)

A point $\mathbf{x} \in \mathcal{X}$ is a spurious stationary point if and only if

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

- We have demonstrated that spurious stationary points are ubiquitous.
- Can we provide a satisfactory stationarity measure? -> New Bregman-type algorithms...

He Chen*, Jiajin Li*, Anthony Man-Cho So*. Spurious Stationarity and Hardness Results for Mirror Descent <http://arxiv.org/abs/2404.08073>

Proposition (Characterization of spurious stationary points)

A point $\mathbf{x} \in \mathcal{X}$ is a spurious stationary point if and only if

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

- We have demonstrated that spurious stationary points are ubiquitous.
- Can we provide a satisfactory stationarity measure? -> New Bregman-type algorithms...

He Chen*, Jiajin Li*, Anthony Man-Cho So*. Spurious Stationarity and Hardness Results for Mirror Descent <http://arxiv.org/abs/2404.08073>

Proposition (Characterization of spurious stationary points)

A point $\mathbf{x} \in \mathcal{X}$ is a spurious stationary point if and only if

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

- We have demonstrated that spurious stationary points are ubiquitous.
- Can we provide a satisfactory stationarity measure? -> New Bregman-type algorithms...

He Chen*, Jiajin Li*, Anthony Man-Cho So*. Spurious Stationarity and Hardness Results for Mirror Descent <http://arxiv.org/abs/2404.08073>

Proposition (Characterization of spurious stationary points)

A point $\mathbf{x} \in \mathcal{X}$ is a spurious stationary point if and only if

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

- We have demonstrated that spurious stationary points are ubiquitous.
- Can we provide a satisfactory stationarity measure? -> New Bregman-type algorithms...

He Chen*, Jiajin Li*, Anthony Man-Cho So*. Spurious Stationarity and Hardness Results for Mirror Descent <http://arxiv.org/abs/2404.08073>

Thank you for your listening!
Any questions?