# Accelerated Bregman proximal gradient methods for relatively smooth convex optimization

Filip Hanzely [*†]    Peter Richtárik [*‡§]    Lin Xiao [¶]

## Abstract

We consider the problem of minimizing the sum of two convex functions: one is differentiable and relatively smooth with respect to a reference convex function, and the other can be nondifferentiable but simple to optimize. We investigate a triangle scaling property of the Bregman distance generated by the reference convex function and present accelerated Bregman proximal gradient (ABPG) methods that attain an $O(k^{-\gamma})$ convergence rate, where $\gamma \in (0, 2]$ is the *triangle scaling exponent* (TSE) of the Bregman distance. For the Euclidean distance, we have $\gamma = 2$ and recover the convergence rate of Nesterov's accelerated gradient methods. For non-Euclidean Bregman distances, the TSE can be much smaller (say $\gamma \leq 1$), but we show that a relaxed definition of *intrinsic* TSE is always equal to 2. We exploit the intrinsic TSE to develop adaptive ABPG methods that converge much faster in practice. Although theoretical guarantees on a fast convergence rate seem to be out of reach in general, our methods obtain empirical $O(k^{-2})$ rates in numerical experiments on several applications and provide posterior numerical certificates for the fast rates.

**Keywords:** convex optimization, relative smoothness, Bregman divergence, proximal gradient methods, accelerated gradient methods.

## 1    Introduction

Let $\mathbb{R}^n$ be the $n$-dimensional real Euclidean space endowed with inner product $\langle x, y \rangle = \sum_{i=1}^{n} x^{(i)} y^{(i)}$ and the Euclidean norm $\|x\| = \sqrt{\langle x, x \rangle}$. We consider optimization problems of the form

$$\underset{x \in C}{\text{minimize}} \ \big\{ F(x) := f(x) + \Psi(x) \big\}, \tag{1}$$

where $C$ is a closed convex set in $\mathbb{R}^n$, and $f$ and $\Psi$ are proper, closed convex functions. We assume that $f$ is differentiable on an open set that contains the relative interior of $C$ (denoted as $\text{rint}\, C$). For the development of first-order methods, we also assume that $C$ and $\Psi$ are *simple*, whose precise meaning will be explained in the context of specific algorithms.

---

[*]Division of Computer, Electrical and Mathematical Sciences and Engineering (CEMSE), King Abdullah University of Science and Technology (KAUST), Thuwal, Kingdom of Saudi Arabia. Emails: `filip.hanzely@kaust.edu.sa`, `peter.richtarik@kaust.edu.sa`

[†]Currently at Toyota Technological Institute at Chicago (TTIC). Email: `filip@ttic.edu`

[‡]School of Mathematics, The University of Edinburgh, Edinburgh, United Kingdom.

[§]Moscow Institute of Physics and Technology, Dolgoprudny, Russia.

[¶]Work done while at Microsoft Research, Redmond, Washington, United States. Email: `lin.xiao@gmail.com`

First-order methods for solving (1) are often based on the idea of minimizing a simple approximation of the objective $F$ during each iteration. Specifically, in the *proximal gradient method*, we start with an initial point $x_0 \in \operatorname{rint} C$ and generate a sequence $x_k$ for $k = 1, 2, \ldots$ with

$$x_{k+1} = \underset{x \in C}{\arg\min} \left\{ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L_k}{2} \|x - x_k\|^2 + \Psi(x) \right\}, \tag{2}$$

where $L_k > 0$ for all $k \geq 0$. Here, we use the gradient $\nabla f(x_k)$ to construct a local quadratic approximation of $f$ around $x_k$ while leaving $\Psi$ untouched. Our assumption that $C$ and $\Psi$ are simple means that the minimization problem in (2) can be solved efficiently, especially if it admits a closed-form solution.

Assuming that $F$ is bounded below, convergence of the proximal gradient method can be established if $F(x_{k+1}) \leq F(x_k)$ for all $k \in \mathbb{N}$. A sufficient condition for this to hold is that the quadratic approximation of $f$ in (2) is an upper approximation (majorization). This is the basic idea behind many general methods for nonlinear optimization. To this end, a common assumption is for the gradient of $f$ to satisfy a uniform Lipschitz condition, i.e., there exists a constant $L_f$ such that

$$\|\nabla f(x) - \nabla f(y)\| \leq L_f \|x - y\|, \quad \forall\, x, y \in \operatorname{rint} C. \tag{3}$$

This smoothness assumption implies (see, e.g., [25, Lemma 1.2.3])

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L_f}{2} \|x - y\|^2, \quad \forall\, x \in C, \ y \in \operatorname{rint} C. \tag{4}$$

Therefore, setting $L_k = L_f$ for all $k \in \mathbb{N}$ ensures that the quadratic approximation of $f$ in (2) is always an upper bound of $f$, which implies $F(x_{k+1}) \leq F(x_k)$ for all $k \in \mathbb{N}$. Moreover, it can be shown that the proximal gradient method enjoys an $O(k^{-1})$ convergence rate, i.e.,

$$F(x_k) - F(x) \leq \frac{L_f}{k} \frac{\|x - x_0\|^2}{2}, \quad \forall\, x \in C. \tag{5}$$

See, e.g., [6], [27] and [5, Chapter 10]. Under the same assumption, accelerated proximal gradient methods ([23, 25, 2, 6, 33, 27]) can achieve a faster $O(k^{-2})$ convergence rate:

$$F(x_k) - F(x) \leq \frac{4L_f}{(k+2)^2} \frac{\|x - x_0\|^2}{2}, \quad \forall\, x \in C, \tag{6}$$

which is optimal (up to a constant factor) for this class of convex optimization problems [22, 25].

## 1.1 Optimization of relatively smooth functions

While the uniform smoothness condition (3) is central in the development and analysis of first-order methods, there are many applications where the objective function does not have this property, despite being convex and differentiable. For example, in D-optimal experiment design (e.g., [19, 1]) and Poisson inverse problems (e.g., [14, 7]), the objective functions involve the logarithm in the form of log-determinant or relative entropy, whose gradients may blow up towards the boundary of the feasible region. In order to develop efficient first-order algorithms for solving such problems, the notion of *relative smoothness* was introduced by several recent works [3, 21, 34].

Let $h$ be a strictly convex function that is differentiable on an open set containing $\operatorname{rint} C$. The Bregman distance associated with $h$, originated in [9] and popularized by [10, 11], is defined as

$$D_h(x, y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle, \quad \forall\, x \in \operatorname{dom} h, \ y \in \operatorname{rint} \operatorname{dom} h.$$

**Definition 1.** *The function $f$ is called $L$-smooth relative to $h$ on $C$ if there is an $L > 0$ such that*

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + LD_h(x, y), \quad \forall\, x \in C,\ y \in \operatorname{rint} C. \tag{7}$$

As shown in [3] and [21], this notion of relative smoothness is equivalent to the following statements:

- $Lh - f$ is a convex function on $C$.
- If both $f$ and $h$ are twice differentiable, then $\nabla^2 f(x) \preceq L\nabla^2 h(x)$ for all $x \in \operatorname{rint} C$.
- $\langle \nabla f(x) - \nabla f(y), x - y \rangle \leq L\langle \nabla h(x) - \nabla h(y), x - y \rangle$ for all $x, y \in \operatorname{rint} C$.

The definition of relative smoothness in (7) gives an upper approximation of $f$ that is similar to (4). In fact, (4) is a special case of (7) with $h = (1/2)\|x\|^2$ and $D_h(x, y) = (1/2)\|x - y\|^2$. Therefore it is natural to consider a more general algorithm by replacing the squared Euclidean distance in (2) with a Bregman distance:

$$x_{k+1} = \arg\min_{x \in C} \left\{ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + L_k D_h(x, x_k) + \Psi(x) \right\}. \tag{8}$$

Here, our assumption that $C$ and $\Psi$ are simple means that the minimization problem in (8) can be solved efficiently. Similar to the proximal gradient method (2), this algorithm can also be interpreted through operator splitting mechanism: it is the composition of a Bregman proximal step and a Bregman gradient step (see details in [3, Section 3.1]). Therefore, it is called the *Bregman proximal gradient* (BPG) method [32].

Under the relative smoothness condition (7), setting $L_k = L$ ensures that the function being minimized in (8) is a majorization of $F$, which implies $F(x_{k+1}) \leq F(x_k)$ for all $k \in \mathbb{N}$. It was first shown in [8] (for the case $\Psi \equiv 0$) that the BGD method has a $O(k^{-1})$ convergence rate:

$$F(x_k) - F(x) \leq \frac{L}{k} D_h(x, x_0), \quad \forall\, x \in \operatorname{dom} h.$$

This is a generalization of (5). The same convergence rate for the general case (with nontrivial $\Psi$) is obtained in [3], where the authors also discussed the effect of a symmetry measure for the Bregman distance. Similar results are also obtained in [21] and [34]. In addition, [21] introduced the notion of relative strong convexity and obtained linear convergence of the BPG method when both relative smoothness and relative strong convexity hold. More recently, [17] studied stochastic gradient descent and randomized coordinate descent methods in the relatively smooth setting, and [20] extended this framework to minimize relatively continuous convex functions.

A natural question is whether the $O(k^{-1})$ rate can be improved with first-order methods under the relative smoothness assumption, especially whether the accelerated $O(k^{-2})$ rate can be achieved [21, 32]. Very recently, it is shown by Dragomir et al. [15] that the $O(k^{-1})$ rate is optimal for the class of relatively smooth functions, thus cannot be improved in general. However, we note that the class of relatively smooth functions is very broad, containing differentiable functions whose gradients has arbitrarily large Lipschitz constants. Indeed, the worst-case function constructed in [15] to prove the lower bound is obtained by smoothing a nonsmooth function, which demonstrate pathological nonsmooth behavior. This is in sharp contrast to the situation under the uniform Lipschitz assumption, which uses a *fixed* quadratic function as the relatively smooth measure.

Ideally, it would be most informative to derive both upper and lower bounds on the convergence rate of first-order methods for every *fixed* function $h$ in the relatively smooth setting, or at least for the popular ones that are frequently encountered in application (such as the KL divergence). It is plausible that the achievable convergence rates for particular functions $h$ (more likely particular combinations of $f$ and $h$) can be better than $O(k^{-1})$ in theory or at least in practice. A full spectrum investigation is beyond the scope of this paper. Instead, we study a structural property of general Bregman divergences called *triangle scaling* and develop adaptive first-order methods that, although without a priori guarantee, often demonstrate the $O(k^{-2})$ convergence rate empirically in many applications. Moreover, these methods produce simple *numerical certificates* of the fast rates whenever they happen.

## 1.2  Contributions and outline

First, in Section 2, we study a triangle-scaling property for general Bregman distances and define the triangle-scaling exponent (TSE) $\gamma > 0$, which is key in characterizing the convergence rates of first-order methods in the relatively smooth setting. We estimate the value of $\gamma$ for several Bregman distances that appear frequently in applications. Moreover, we define an intrinsic triangle-scaling exponent $\gamma_{\mathrm{in}}$ and show that $\gamma_{\mathrm{in}} = 2$ for all $h$ that is twice continuously differentiable.

In Section 3, we propose a basic accelerated Bregman proximal gradient (ABPG) method that attains an $O(k^{-\gamma})$ convergence rate, where $\gamma \leq 2$ is the TSE of the Bregman divergence. More specifically, under the assumption (7), the basic ABPG method produces a sequence $\{x_k\}$ satisfying

$$F(x_k) - F(x) \leq \left(\frac{\gamma}{k+\gamma}\right)^{\gamma} L D_h(x, x_0), \quad \forall\, x \in \mathrm{dom}\, h. \tag{9}$$

The exact value of $\gamma$ depends on a *triangle scaling property* of the Bregman distance. For $D_h(x, y) = (1/2)\|x - y\|^2$, we have $\gamma = 2$ and $L = L_f$, hence the result in (9) recovers that in (6). We also give an adaptive variant that can automatically search for the largest possible $\gamma$ for which the convergence rate in (9) holds for finite $k$ even though $\gamma$ is larger than the TSE.

In Section 4, we develop an adaptive ABPG method that automatically adjust an additional gain factor in order to work with the intrinsic TSE $\gamma_{\mathrm{in}} = 2$. If the geometric mean of the gains obtained at all the iterations up to $k$ is a small constant, say $O(1)$, then they constitute numerical certificates that the algorithm has enjoyed an empirical $O(k^{-2})$ convergence rate.

In Section 5, we present an accelerated Bregman dual-averaging algorithm that has similar convergence rates as the basic ABPG method, but omit discussions of its adaptive variants.

Finally, in Section 6, we present numerical experiments with three applications: the D-optimal experiment design problem, a Poisson linear inverse problem, and relative-entropy nonnegative regression. In all experiments, the ABPG methods, especially the adaptive variants, demonstrate superior performance compared with the BPG method. Moreover, we obtain numerical certificates for the empirical $O(k^{-2})$ rate in all our experiments.

**Related work.**  The relative smoothness condition directly extends the upper approximation property (4) with more general Bregman distances. Nesterov [28] took an alternative approach by extending the Lipschitz condition (3). Specifically, he considered functions with Hölder continuous gradients with a parameter $\nu \in [0, 1]$:

$$\|\nabla f(x) - \nabla f(y)\|_* \leq L_\nu \|x - y\|^\nu, \quad x, y \in C,$$

and obtained $O(k^{-(1+\nu)/2})$ rate with a universal gradient method and $O(k^{-(1+3\nu)/2})$ rate with accelerated schemes. These methods are called "universal" because they do not assume the knowledge of $\nu$ and automatically ensure the best possible rate of convergence. The accelerated $O(k^{-(1+3\nu)/2})$ rate interpolates between $O(k^{-1/2})$ and $O(k^{-2})$ with $\nu \in [0, 1]$. There seems to be no simple connection or correspondence between the Hölder smoothness property and the combination of relative smoothness and the triangle scaling property studied in this paper.

Gutman and Peña [16] studied iteration complexity of first-order methods using a general framework of perturbed Fenchel duality. Their framework provides alternative derivations of the convergence rates of Bregman proximal gradient methods under the relative smooth setting and the ones under Hölder continuity assumption.

**Technical assumptions.** Development and analysis of optimization methods in the relatively smooth setting require some delicate assumptions in order to cover many interesting applications without loss of rigor. Here we adopt the same assumptions made in [3] regarding problem (1).

**Assumption A.** *Suppose that $C$ is a closed convex set in $\mathbb{R}^n$ and $h : \mathbb{R}^n \to (-\infty, +\infty]$ is strictly convex and differentiable on an open set containing* $\mathrm{rint}\, C$. *Moreover,*

1. *$h$ is of Legendre type [31, Section 26]. In other words, it is essentially smooth and strictly convex in* $\mathrm{rint}\,\mathrm{dom}\, h$. *Essential smoothness means that it is differentiable and $\|\nabla h(x_k)\| \to \infty$ for every sequence $\{x_k\}_{k \in \mathbb{N}}$ converging to a boundary point of* $\mathrm{dom}\, h$.

2. *$f : \mathbb{R}^n \to (-\infty, \infty]$ is a proper and closed convex function, and it is differentiable on* $\mathrm{rint}\, C$.

3. *$\Psi : \mathbb{R}^n \to (-\infty, \infty]$ is a proper and closed convex function, and $\mathrm{dom}\, \Psi \cap \mathrm{rint}\,\mathrm{dom}\, h \neq \emptyset$.*

4. *$\inf_{x \in C}\{f(x) + \Psi(x)\} > -\infty$, i.e., problem (1) is bounded below.*

5. *The BPG step (8) is well posed, meaning that $x_{k+1}$ is unique and belongs to* $\mathrm{rint}\,\mathrm{dom}\, h$.

Sufficient conditions for the well-posedness of (8) are given in [3, Lemma 2]. The same conditions also ensure that our proposed accelerated methods are well-posed.

## 2 Triangle scaling of Bregman distance

In this section, we define the *triangle scaling property* for Bregman distances and discuss two different notions of *triangle scaling exponent* (TSE).

**Definition 2.** *Let $h$ be a convex function that is differentiable on* $\mathrm{rint}\,\mathrm{dom}\, h$. *The Bregman distance $D_h$ has the* triangle scaling property *if there is a constant $\gamma > 0$ such that for all $x, z, \tilde{z} \in \mathrm{rint}\,\mathrm{dom}\, h$,*

$$D_h\big((1-\theta)x + \theta z,\ (1-\theta)x + \theta\tilde{z}\big) \ \leq\ \theta^\gamma D_h(z, \tilde{z}), \qquad \forall\, \theta \in [0, 1]. \tag{10}$$

*We call $\gamma$ a* uniform *triangle scaling exponent (TSE) of $D_h$.*

Figure 1 gives a geometric illustration of the points involved in the above definition.

If $D_h(x, y)$ is jointly convex in $(x, y)$, then the inequality (10) holds with $\gamma = 1$ because

$$D_h\big((1-\theta)x + \theta z,\ (1-\theta)x + \theta\tilde{z}\big) \leq (1-\theta)D_h(x, x) + \theta D_h(z, \tilde{z}) = \theta D_h(z, \tilde{z}).$$
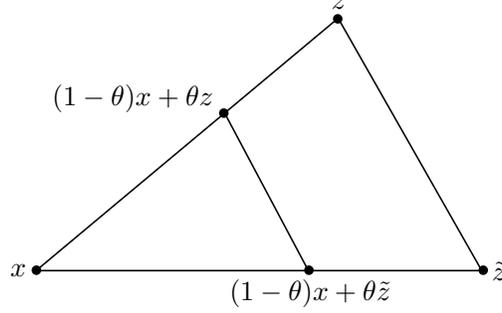
Figure 1: Illustration of different points in the triangle scaling property.

Therefore it is useful to study jointly convex Bregman distances. Suppose $h : \mathbb{R} \to (-\infty, \infty]$ is strictly convex and twice continuously differentiable on an open interval in $\mathbb{R}$. Let $h''$ denotes the second derivative of $h$. It was shown in [4] that the Bregman distance $D_h(\cdot, \cdot)$ is jointly convex if and only if $1/h''$ is concave. This result applies directly to separable functions which can be written as $h(x) = \sum_{i=1}^{n} h_i(x^{(i)})$. If $1/h_i''$ is concave for each $i = 1, \dots, n$, then we conclude that $D_h$ has a uniform TSE of at least 1. Below are some specific examples:

- *The squared Euclidean distance.* Let $h(x) = (1/2)\|x\|_2^2$ and $D_h(x, y) = (1/2)\|x - y\|_2^2$. Obviously, here $D_h$ is jointly convex in its two arguments. But it is also easy to see that

$$\frac{1}{2}\left\| (1 - \theta)x + \theta z - \left((1 - \theta)x + \theta \tilde{z}\right)\right\|_2^2 = \frac{1}{2}\|\theta(z - \tilde{z})\|_2^2 = \theta^2 \frac{1}{2}\|z - \tilde{z}\|_2^2.$$

  Therefore the squared Euclidean distance has a uniform TSE $\gamma = 2$, which is much larger than 1 obtained by following the jointly convex argument.

- *Bregman distance induced by strongly convex and smooth functions.* If $h$ is $\mu$-strongly convex and $L$-smooth over its domain, then the inequality (10) would hold with $\gamma = 2$ if the right-hand side is multiplied by an additional factor $G = L/\mu$, which is the condition number of $h$. We will prove this fact in Section 2.2.

- *The generalized Kullback-Leibler (KL) divergence.* Let $h$ be the negative Boltzmann-Shannon entropy: $h(x) = \sum_{i=1}^{n} x^{(i)} \log x^{(i)}$ defined on $\mathbb{R}_+^n$. The Bregman distance associated with $h$ is

$$D_{\mathrm{KL}}(x, y) = \sum_{i=1}^{n} \left( x^{(i)} \log\left(\frac{x^{(i)}}{y^{(i)}}\right) - x^{(i)} + y^{(i)} \right). \tag{11}$$

  Since $1/h_i'' = x^{(i)}$ is linear thus concave for each $i$, we conclude that $D_{\mathrm{KL}}(x, y)$ is jointly convex in $(x, y)$, which implies that it has a uniform TSE $\gamma = 1$.

- *The Itakura-Saito (IS) distance.* The IS distance is the Bregman distance generated by Burg's entropy $h(x) = \sum_{i=1}^{n} -\log(x^{(i)})$ with dom $h = \mathbb{R}_{++}^n$:

$$D_{\mathrm{IS}}(x, y) = \sum_{i=1}^{n} \left( -\log\left(\frac{x^{(i)}}{y^{(i)}}\right) + \frac{x^{(i)}}{y^{(i)}} - 1 \right). \tag{12}$$

6

Since $1/h_i'' = (x^{(i)})^2$ is not concave, we conclude that $D_{\mathrm{IS}}(\cdot, \cdot)$ is not jointly convex. Hence if it has a uniform TSE, then it is likely to be less than 1. In fact, it can be easily checked numerically that any $\gamma > 0.5$ is not a uniform TSE for $D_{\mathrm{IS}}$ when $G = 1$.

- *Bregman distance based on polynomial kernels.* Reference functions of the form $h(x) = (1/p)\|x\|^p$ for some $p \geq 2$ recently attracted lots of attention following Nesterov's work on tensor methods in convex optimization [29]. In general, the global TSEs for the induced Bregman divergence can be less than 1 for $p > 2$. However, the modified reference function $h(x) = (1/2)\|x\|^2 + (1/p)\|x\|^p$ for $p \geq 4$ has $\gamma > 1$, or $\gamma = 2$ with an additional factor on the right-hand side of (10), over a bounded domain. We will give detailed analysis for the case $p = 4$ in Section 2.2, after introducing a relaxed version of TSE.

We observe that the largest uniform TSEs are quite different for the Bregman distances listed above. An important question is: Are these differences essential such that they lead to different convergence rates if different Bregman distances are used in an accelerated algorithm? It would be ideal to derive an intrinsic characterization that is common for most Bregman distances and essential for convergence analysis of accelerated algorithms.

## 2.1 The intrinsic triangle-scaling exponent

For any fixed triple $\{x, z, \tilde{z}\} \subset \operatorname{rint} \operatorname{dom} h$, we consider a relaxed version of triangle scaling:

$$D_h\big((1-\theta)x + \theta z, \ (1-\theta)x + \theta\tilde{z}\big) \ \leq \ G(x, z, \tilde{z})\, \theta^\gamma D_h(z, \tilde{z}), \qquad \forall\, \theta \in [0, 1], \tag{13}$$

where $G(x, z, \tilde{z})$ depends on the triple $\{x, z, \tilde{z}\}$ but does not depend on $\theta$. The intrinsic TSE of $D_h$, denoted $\gamma_{\mathrm{in}}$, is the largest $\gamma$ such that (13) holds with some finite $G(x, z, \tilde{z})$ for all triples $\{x, z, \tilde{z}\} \subset \operatorname{rint} \operatorname{dom} h$.

Notice that when $\theta$ is bounded away from 0, we can always find sufficiently large $G(x, z, \tilde{z})$ to make the inequality in (13) hold with any value of $\gamma$. Therefore, the intrinsic TSE is determined only by the asymptotic behavior of $D_h\big((1-\theta)x + \theta z, (1-\theta)x + \theta\tilde{z}\big)$ when $\theta \to 0$. A more precise definition is as follows.

**Definition 3.** *The intrinsic TSE of $D_h$, denoted $\gamma_{\mathrm{in}}$, is the largest $\gamma$ such that for all $x, z, \tilde{z} \in \operatorname{rint} \operatorname{dom} h$,*

$$\limsup_{\theta \to 0} \frac{D_h\big((1-\theta)x + \theta z, (1-\theta)x + \theta\tilde{z}\big)}{\theta^\gamma} \ < \ \infty.$$

We show that a broad family of Bregman distances share the same intrinsic TSE $\gamma_{\mathrm{in}} = 2$.

**Proposition 1.** *If $h$ is convex and twice continuously differentiable on $\operatorname{rint} \operatorname{dom} h$, then the intrinsic TSE of the Bregman distance $D_h$ is 2. Specifically, for any $\{x, z, \tilde{z}\} \subset \operatorname{rint} \operatorname{dom} h$, we have*

$$\lim_{\theta \to 0} \frac{D_h\big((1-\theta)x + \theta z, \ (1-\theta)x + \theta\tilde{z}\big)}{\theta^2} = \frac{1}{2}\langle \nabla^2 h(x)(z - \tilde{z}), z - \tilde{z}\rangle. \tag{14}$$

*Proof.* Consider the limit in (14), since both the numerator $D_h\big((1-\theta)x + \theta z, \ (1-\theta)x + \theta\tilde{z}\big)$ and the denominator $\theta^2$ converge to zero as $\theta \to 0$, we apply L'Hospital's rule. First, by definition of

7

the Bregman distance, we have

$$
\begin{aligned}
&D_h\big((1-\theta)x + \theta z,\ (1-\theta)x + \theta\tilde{z}\big)\\
&= D_h\big(x + \theta(z-x),\ x + \theta(\tilde{z}-x)\big)\\
&= h\big(x + \theta(z-x)\big) - h\big(x + \theta(\tilde{z}-x)\big) - \big\langle \nabla h\big(x + \theta(\tilde{z}-x)\big), \theta(z-\tilde{z})\big\rangle.
\end{aligned}
$$

The derivative of $D_h\big((1-\theta)x + \theta z,\ (1-\theta)x + \theta\tilde{z}\big)$ with respect to $\theta$ is

$$
\frac{d}{d\theta}D_h\big((1-\theta)x + \theta z,\ (1-\theta)x + \theta\tilde{z}\big) = A(\theta) - \big\langle \nabla^2 h(x + \theta(\tilde{z}-x))(\tilde{z}-x), \theta(z-\tilde{z})\big\rangle,
$$

where

$$
A(\theta) = \big\langle \nabla h(x + \theta(z-x)), z-x\big\rangle - \big\langle \nabla h(x + \theta(\tilde{z}-x)), \tilde{z}-x\big\rangle - \big\langle \nabla h(x + \theta(\tilde{z}-x)), z-\tilde{z}\big\rangle.
$$

Therefore,

$$
\begin{aligned}
\lim_{\theta\to 0} \frac{D_h\big((1-\theta)x + \theta z, (1-\theta)x + \theta\tilde{z}\big)}{\theta^2} &= \lim_{\theta\to 0} \frac{A(\theta) - \big\langle \nabla^2 h\big(x + \theta(\tilde{z}-x)\big)(\tilde{z}-x), \theta(z-\tilde{z})\big\rangle}{2\theta}\\
&= \lim_{\theta\to 0}\frac{A(\theta)}{2\theta} - \lim_{\theta\to 0}\frac{\big\langle \nabla^2 h\big(x + \theta(\tilde{z}-x)\big)(\tilde{z}-x), z-\tilde{z}\big\rangle}{2}\\
&= \lim_{\theta\to 0}\frac{A(\theta)}{2\theta} - \frac{1}{2}\big\langle \nabla^2 h(x)(\tilde{z}-x), z-\tilde{z}\big\rangle. \qquad (15)
\end{aligned}
$$

Notice that

$$
\lim_{\theta\to 0} A(\theta) = \big\langle \nabla h(x), z-x\big\rangle - \big\langle \nabla h(x), \tilde{z}-x\big\rangle - \big\langle \nabla h(x), z-\tilde{z}\big\rangle = 0,
$$

so we apply L'Hospital's rule again:

$$
\lim_{\theta\to 0}\frac{A(\theta)}{2\theta} = \frac{\big\langle \nabla^2 h(x)(z-x), z-x\big\rangle - \big\langle \nabla^2 h(x)(\tilde{z}-x), \tilde{z}-x\big\rangle - \big\langle \nabla^2 h(x)(\tilde{z}-x), z-\tilde{z}\big\rangle}{2}.
$$

Plugging the last equality into (15) and after some simple algebra, we arrive at (14). $\qquad\square$

According to Proposition 1, the examples we considered earlier, including the generalized KL-divergence and the IS-distance, share the same intrinsic TSE $\gamma_{\text{in}} = 2$. Proposition 1 also implies that the largest uniform TSE cannot exceed 2.

## 2.2  Bounding the triangle-scaling gain

In our analysis of accelerated algorithms in the relatively smooth setting, it is crucial to bound the triangle scaling gain $G(x, z, \tilde{z})$. Here we derive a general bound based on the relative scaling of the Hessians of $h$ at different points. First, by the second-order Taylor expansion (mean value theorem), we have

$$
D_h(x, y) = h(x) - h(y) - \langle \nabla h(y), x-y\rangle = \frac{1}{2}(x-y)^T \nabla^2 h(w)(x-y),
$$

where $w = x + t(y-x)$ for some $t \in [0, 1]$, which we denote as $w \in [x, y]$. Consequently, if we define

$$
G_\theta(x, z, \tilde{z}) := \frac{D_h\big((1-\theta)x + \theta z, (1-\theta)x + \theta\tilde{z}\big)}{\theta^2 \cdot D_h(z, \tilde{z})},
$$

8

then for some $u \in [(1 - \theta)x + \theta z, (1 - \theta)x + \theta \tilde{z}]$ and $v \in [z, \tilde{z}]$, we have

$$G_\theta(x, z, \tilde{z}) = \frac{\frac{1}{2}\theta^2(z - \tilde{z})\nabla^2 h(u)(z - \tilde{z})}{\theta^2 \cdot \frac{1}{2}(z - \tilde{z})^T \nabla^2 h(v)(z - \tilde{z})} = \frac{(z - \tilde{z})\nabla^2 h(u)(z - \tilde{z})}{(z - \tilde{z})^T \nabla^2 h(v)(z - \tilde{z})}.$$

The last expression does not depend on $\theta$ explicitly, but through $u \in [(1 - \theta)x + \theta z, (1 - \theta)x + \theta \tilde{z}]$. If $\theta \to 0$, then we have $u \to x$.

In general, let's assume $u, v \in \operatorname{rint} \operatorname{dom} h$ and $\nabla^2 h(v)$ is non-singular. Then we have

$$G_\theta(x, z, \tilde{z}) \leq \lambda_{\max}\left(\nabla^2 h(v)^{-1/2}\nabla^2 h(u)\nabla^2 h(v)^{-1/2}\right), \tag{16}$$

where $\lambda_{\max}(\cdot)$ denotes the maximum eigenvalue of a positive semidefinite matrix. Therefore, the triangle-scaling gain is bounded by how close the two Hessians $\nabla^2 h(u)$ and $\nabla^2 h(v)$ are. Since convex quadratic functions of the form $h(x) = (1/2)x^T A x + b^T x + c$ has constant Hessian $\nabla^2 h(x) = A$, the triangle scaling gain is always 1, independent of $\theta$.

More generally, if $h$ is strongly convex and smooth (the eigenvalues of its Hessian have positive lower and upper bounds), then the triangle scaling gain is bounded by its condition number, i.e., the ratio between the upper and lower bounds on the Hessian eigenvalues. Otherwise, the gain can be unbounded without any proximity assumption on the three points $(x, z, \tilde{z})$.

Next we consider the polynomial reference function $h(x) = (1/4)\|x\|^4$, which does not have bounded Hessian. In this case, we have $\nabla h(x) = \|x\|^2 x$ and $\nabla^2 h(x) = \|x\|^2 \cdot I + 2xx^T$, where $I$ is the identity matrix. Clearly $\|x\|^2 \cdot I \preceq \nabla^2 h(x) \preceq 3\|x\|^2 \cdot I$. According to (16), we have

$$G_\theta(x, z, \tilde{z}) \leq 3\frac{\|u\|^2}{\|v\|^2},$$

for some $u \in [(1 - \theta)x + \theta z, (1 - \theta)x + \theta \tilde{z}]$ and $v \in [z, \tilde{z}]$. Therefore, it is not hard to construct examples with $v \approx 0$ thus the triangle scaling gain can be unbounded, even if the points $(x, z, \tilde{z})$ are close in a small neighborhood (near the origin).

As a simple fix, we consider $h(x) = (1/2)\|x\|^2 + (1/4)\|x\|^4$, whose Hessian is $\nabla^2 h(x) = (1 + \|x\|^2)I + 2xx^T$ and it satisfies $(1 + \|x\|^2)I \preceq \nabla^2 h(x) \preceq (1 + 3\|x\|^2)I$. Therefore, according to (16),

$$G_\theta(x, z, \tilde{z}) \leq \frac{1 + 3\|u\|^2}{1 + \|v\|^2}.$$

In this case, it is clear that $G_\theta(x, z, \tilde{z})$ (associated with $\gamma_{\text{in}} = 2$) is always bounded if the three points $(x, z, \tilde{z})$ are bounded, even if $\|v\| = 0$. In particular, if the domain of consideration, $\operatorname{dom}\Psi$, is bounded with radius $R$ from the origin, then $G_\theta(x, z, \tilde{z}) \leq 1 + 3R^2$.

## 3 Accelerated Bregman proximal gradient method

In this section, we present an accelerated Bregman proximal gradient (ABPG) method for solving problem (1), and analyze its convergence rate under the uniform triangle-scaling property. Adaptive variants based on the intrinsic TSE are developed in Section 4.

To simplify notation, we define a lower approximation of $F(x) = f(x) + \Psi(x)$ by linearizing $f$ at a given point $y$:

$$\ell(x|y) := f(y) + \langle \nabla f(y), x - y \rangle + \Psi(x).$$

---

**Algorithm 1:** Accelerated Bregman proximal gradient (ABPG) method

---

**input:** initial point $x_0 \in \operatorname{rint} C$ and $\gamma \geq 1$.

initialize: $z_0 = x_0$ and $\theta_0 = 1$.

**for** $k = 0, 1, 2, \ldots$ **do**

1     $y_k = (1 - \theta_k)x_k + \theta_k z_k$

2     $z_{k+1} = \arg\min_{z \in C}\{\ell(z|y_k) + \theta_k^{\gamma-1}LD_h(z, z_k)\}$

3     $x_{k+1} = (1 - \theta_k)x_k + \theta_k z_{k+1}$

4     choose $\theta_{k+1} \in (0, 1]$ such that $\frac{1-\theta_{k+1}}{\theta_{k+1}^{\gamma}} \leq \frac{1}{\theta_k^{\gamma}}$

**end**

---

If $f$ is $L$-smooth relative to $h$ (Definition 1), then we have both a lower and an upper approximation:

$$\ell(x|y) \;\leq\; F(x) \;\leq\; \ell(x|y) + LD_h(x, y). \tag{17}$$

Algorithm 1 describes the ABPG method. Its input parameters include a uniform TSE $\gamma$ of $D_h$ and an initial point $x_0 \in \operatorname{rint} C$. The sequence $\{\theta_k\}_{k \in \mathbb{N}}$ in Algorithm 1 satisfies $0 < \theta_k \leq 1$ and

$$\frac{1 - \theta_{k+1}}{\theta_{k+1}^{\gamma}} \leq \frac{1}{\theta_k^{\gamma}}, \qquad \forall\, k \geq 0. \tag{18}$$

When $\gamma = 2$ and $\Psi \equiv 0$, Algorithm 1 reduces to the IGA (improved interior gradient algorithm) method in [2], which is an extension of Nesterov's accelerated gradient method in [24] to the Bregman proximal setting. It was shown in [2] that the IGA method attains $O(k^{-2})$ rate of convergence under the uniform Lipschitz condition (3). In this paper, we consider the general case $\gamma \in [1, 2]$ under the much weaker relatively smooth condition.

Using the definition of $\ell(\cdot|\cdot)$, line 2 in Algorithm 1 can be written as

$$z_{k+1} = \underset{x \in C}{\arg\min}\left\{f(y_k) + \langle \nabla f(y_k), x - y_k \rangle + \theta_k^{\gamma-1}LD_h(x, z_k) + \Psi(x)\right\}, \tag{19}$$

which is very similar to the BPG step (8). Here the function $f$ is linearized around $y_k$ but the Bregman distance is measured from a different point $z_k$. Therefore it does not fit into the framework of majorization and the sequence $F(x_k)$ may not be monotone decreasing. However, the upper bound in (17) is still crucial to ensure convergence of the algorithm. Under the same assumption that the BPG step is well-posed (Assumption A.5), the ABPG method is also well-posed, meaning that $z_{k+1} \in \operatorname{rint} C$ always and it is unique.

## 3.1 Convergence analysis of ABPG

We show that the ABPG method converges with a sublinear rate of $O(k^{-\gamma})$. First, we state a basic property of optimization with Bregman distance [13, Lemma 3.2].

**Lemma 1.** *For any closed convex function $\varphi : \mathbb{R}^n \to (-\infty, \infty]$ and any $z \in \operatorname{rint} \operatorname{dom} h$, if*

$$z_+ = \underset{x \in C}{\arg\min}\left\{\varphi(x) + D_h(x, z)\right\}$$

10

*and h is differentiable at $z_+$, then*

$$\varphi(x) + D_h(x, z) \geq \varphi(z_+) + D_h(z_+, z) + D_h(x, z_+), \quad \forall x \in \operatorname{dom} h.$$

The following lemma establishes a relationship between the two consecutive steps of Algorithm 1. It is an extension of Proposition 1 in [33] , which uses $\gamma = 2$ under the assumption (3).

**Lemma 2.** *Suppose Assumption A holds, $f$ is L-smooth relative to $h$ on $C$, and $\gamma$ is a uniform TSE of $D_h$. For any $x \in \operatorname{dom} h$, the sequences generated by Algorithm 1 satisfy, for all $k \geq 0$,*

$$\frac{1 - \theta_{k+1}}{\theta_{k+1}^{\gamma}} \big( F(x_{k+1}) - F(x) \big) + LD_h(x, z_{k+1}) \; \leq \; \frac{1 - \theta_k}{\theta_k^{\gamma}} \big( F(x_k) - F(x) \big) + LD_h(x, z_k). \tag{20}$$

*Proof.* First, using the upper approximation in (17) and line 1 and line 3 in Algorithm 1, we have

$$
\begin{aligned}
F(x_{k+1}) &\leq \ell(x_{k+1}|y_k) + LD_h(x_{k+1}, y_k) \\
&= \ell(x_{k+1}|y_k) + LD_h\big((1 - \theta_k)x_k + \theta_k z_{k+1}, (1 - \theta_k)x_k + \theta_k z_k\big) \\
&\leq \ell(x_{k+1}|y_k) + \theta_k^{\gamma} LD_h(z_{k+1}, z_k),
\end{aligned}
\tag{21}
$$

where in the last inequality we used the triangle-scaling property (10). Using $x_{k+1} = (1 - \theta_k)x_k + \theta_k z_{k+1}$ and convexity of $\ell(\cdot|y_k)$, we have

$$
\begin{aligned}
F(x_{k+1}) &\leq (1 - \theta_k)\ell(x_k|y_k) + \theta_k \ell(z_{k+1}|y_k) + \theta_k^{\gamma} LD_h(z_{k+1}, z_k) \\
&= (1 - \theta_k)\ell(x_k|y_k) + \theta_k \left( \ell(z_{k+1}|y_k) + \theta_k^{\gamma-1} LD_h(z_{k+1}, z_k) \right).
\end{aligned}
\tag{22}
$$

Now applying Lemma 1 with $\varphi(x) = \ell(x|y_k)/(\theta^{\gamma-1}L)$ yields, for any $x \in \operatorname{dom} h$,

$$\ell(z_{k+1}|y_k) + \theta_k^{\gamma-1} LD_h(z_{k+1}, z_k) \; \leq \; \ell(x|y_k) + \theta_k^{\gamma-1} LD_h(x, z_k) - \theta_k^{\gamma-1} LD_h(x, z_{k+1}).$$

Hence

$$
\begin{aligned}
F(x_{k+1}) &\leq (1 - \theta_k)\ell(x_k|y_k) + \theta_k \left( \ell(x|y_k) + \theta_k^{\gamma-1} LD_h(x, z_k) - \theta_k^{\gamma-1} LD_h(x, z_{k+1}) \right) \\
&= (1 - \theta_k)\ell(x_k|y_k) + \theta_k \ell(x|y_k) + \theta_k^{\gamma} \big( LD_h(x, z_k) - LD_h(x, z_{k+1}) \big) \\
&\leq (1 - \theta_k)F(x_k) + \theta_k F(x) + \theta_k^{\gamma} \big( LD_h(x, z_k) - LD_h(x, z_{k+1}) \big),
\end{aligned}
$$

where in the last inequality we used the lower bound in (17). Subtracting $F(x)$ from both sides of the inequality above, we obtain

$$F(x_{k+1}) - F(x) \; \leq \; (1 - \theta_k)\big( F(x_k) - F(x) \big) + \theta_k^{\gamma} \big( LD_h(x, z_k) - LD_h(x, z_{k+1}) \big).$$

Dividing both sides by $\theta_k^{\gamma}$ and rearranging terms yield

$$\frac{1}{\theta_k^{\gamma}} \big( F(x_{k+1}) - F(x) \big) + LD_h(x, z_{k+1}) \; \leq \; \frac{1 - \theta_k}{\theta_k^{\gamma}} \big( F(x_k) - F(x) \big) + LD_h(x, z_k). \tag{23}$$

Finally applying the condition (18) gives the desired result. □

**Lemma 3.** *The sequence $\theta_k = \frac{\gamma}{k+\gamma}$ for $k = 0, 1, 2, \ldots$ satisfies the condition (18).*

*Proof.* With $\theta_k = \frac{\gamma}{k+\gamma}$, we have

$$\frac{1 - \theta_{k+1}}{\theta_{k+1}^\gamma} = \left(1 - \frac{\gamma}{k+1+\gamma}\right)\left(\frac{k+1+\gamma}{\gamma}\right)^\gamma = \frac{(k+1)(k+1+\gamma)^{\gamma-1}}{\gamma^\gamma} \tag{24}$$

and

$$\frac{1}{\theta_k^\gamma} = \left(\frac{k+\gamma}{\gamma}\right)^\gamma = \frac{(k+\gamma)^\gamma}{\gamma^\gamma}. \tag{25}$$

Recall the weighted arithmetic mean and geometric mean inequality (see, e.g., [18, Section 2.5].), i.e., for any positive real numbers $a$, $b$, $\alpha$ and $\beta$, it holds that

$$a^\alpha b^\beta \le \left(\frac{\alpha a + \beta b}{\alpha + \beta}\right)^{\alpha+\beta}. \tag{26}$$

Setting $a = k + 1$, $b = k + 1 + \gamma$, $\alpha = 1$ and $\beta = \gamma - 1$, we arrive at

$$(k+1)(k+1+\gamma)^{\gamma-1} \le \left(\frac{k+1+(\gamma-1)(k+1+\gamma)}{1+\gamma-1}\right)^{1+\gamma-1} = (k+\gamma)^\gamma,$$

which, together with (24) and (25), implies the inequality (18). $\qquad\square$

A slightly faster converging sequence $\theta_k$ can be obtained by solving the equality in (18). Since there is no closed-form solution in general, we can find $\theta_{k+1}$ as the root of

$$\theta^\gamma - \theta_k^\gamma(1 - \theta) = 0 \tag{27}$$

numerically, say, using Newton's method with $\theta_k$ as the starting point.

**Lemma 4.** *Let $\theta_0 = 1$ and $\theta_{k+1}$ be the solution to (27) for all $k \ge 0$. Then $\theta_k \le \frac{\gamma}{k+\gamma}$ for all $k \ge 0$.*

*Proof.* Let $\vartheta_k = \frac{\gamma}{k+\gamma}$ and define another sequence $\xi_k$ such that $\xi_0 = 1$ and

$$\frac{1 - \xi_{k+1}}{\xi_{k+1}^\gamma} = \frac{1}{\vartheta_k^\gamma}, \qquad \forall\, k \ge 0. \tag{28}$$

Notice that the function

$$\omega(\theta) := \frac{1 - \theta}{\theta^\gamma}$$

is monotone decreasing in $\theta$. Since $\omega(\vartheta_{k+1}) \le 1/\vartheta_k^\gamma$ by Lemma 3 and $\omega(\xi_{k+1}) = 1/\vartheta_k^\gamma$ by (28), we have $\xi_{k+1} \le \vartheta_{k+1}$ for all $k \ge 0$.

Next we prove $\theta_k \le \vartheta_k$ for all $k \ge 0$ by mathematical induction. This obviously holds for $k = 0$ since $\theta_0 = \vartheta_0 = 1$. Suppose $\theta_k \le \vartheta_k$ holds for some $k \ge 0$. Then using the facts $\omega(\theta_{k+1}) = 1/\theta_k^\gamma$ and $\omega(\xi_{k+1}) = 1/\vartheta_k^\gamma$, we obtain $\omega(\theta_{k+1}) \ge \omega(\xi_{k+1})$. Since $\omega$ is monotone decreasing, we conclude that $\theta_{k+1} \le \xi_{k+1}$. Combining with $\xi_{k+1} \le \vartheta_{k+1}$ obtained above, we have $\theta_{k+1} \le \vartheta_{k+1}$. This completes the induction. $\qquad\square$

**Theorem 1.** *Suppose Assumption A holds, $f$ is $L$-smooth relative to $h$ on $C$, and $\gamma$ is a uniform TSE of $D_h$. If $\theta_k \le \frac{\gamma}{k+\gamma}$ for all $k \ge 0$, then the outputs of Algorithm 1 satisfy, for any $x \in \operatorname{dom} h$,*

$$F(x_{k+1}) - F(x) \le \left(\frac{\gamma}{k+\gamma}\right)^\gamma LD_h(x, x_0), \qquad \forall\, k \ge 0.$$

**Algorithm 2:** ABPG method with exponent adaptation (ABPG-e)

---

**input:** initial point $x_0 \in \operatorname{rint} C$, $\gamma_0 \geq 2$, $\gamma_{\min} \geq 0$, and $\delta > 0$.

initialize: $z_0 = x_0$, $\gamma_{-1} = \gamma_0$, and $\theta_0 = 1$.

**for** $k = 0, 1, 2, \dots$ **do**

$\quad y_k = (1 - \theta_k)x_k + \theta_k z_k$

$\quad$ **repeat** for $t = 0, 1, 2, \dots$

$\quad\quad \gamma_k = \max\{\gamma_{k-1} - \delta t, \ \gamma_{\min}\}$

$\quad\quad z_{k+1} = \arg\min_{z \in C}\big\{\ell(z|y_k) + \theta_k^{\gamma_k - 1} LD_h(z, z_k)\big\}$

$\quad\quad x_{k+1} = (1 - \theta_k)x_k + \theta_k z_{k+1}$

$\quad$ **until** $\ f(x_{k+1}) \leq f(y_k) + \langle \nabla f(y_k), x_{k+1} - y_k \rangle + \theta_k^{\gamma_k} LD_h(z_{k+1}, z_k)$

$\quad$ choose $\theta_{k+1} \in (0, 1]$ such that $\frac{1 - \theta_{k+1}}{\theta_{k+1}^{\gamma_k}} \leq \frac{1}{\theta_k^{\gamma_k}}$

**end**

---

*Proof.* A direct consequence of Lemma 2 is, for any $x \in \operatorname{dom} h$,

$$\frac{1 - \theta_k}{\theta_k^{\gamma}}\big(F(x_k) - F(x)\big) + LD_h(x, z_k) \leq \frac{1 - \theta_0}{\theta_0}\big(F(x_0) - F(x)\big) + LD_h(x, z_0).$$

Combining with (23), we have

$$\frac{1}{\theta_k^{\gamma}}\big(F(x_{k+1}) - F(x)\big) + LD_h(x, z_{k+1}) \leq \frac{1 - \theta_0}{\theta_0}\big(F(x_0) - F(x)\big) + LD_h(x, z_0).$$

Using $D_h(x, z_{k+1}) \geq 0$ and the initializations $\theta_0 = 1$ and $z_0 = x_0$, we obtain

$$\frac{1}{\theta_k^{\gamma}}\big(F(x_{k+1}) - F(x)\big) \leq LD_h(x, z_0),$$

which implies

$$F(x_{k+1}) - F(x) \leq \theta_k^{\gamma} LD_h(x, x_0).$$

It remains to apply the condition $\theta_k \leq \frac{\gamma}{k + \gamma}$. $\qquad\qquad\qquad\qquad\qquad\qquad\square$

## 3.2 ABPG method with exponent adaptation

The best convergence rate of the ABPG method is obtained with the largest uniform TSE for the Bregman distance. Since it is often hard to determine the largest TSE, we present in Algorithm 2 a variant of the ABPG method with automatic exponent adaptation, called the ABPG-e method.

This method starts with a large $\gamma_0 \geq 2$. During each iteration $k$, it reduces $\gamma_k$ by a small amount $\delta > 0$ until some stopping criterion is satisfied. An obvious choice for the stopping criterion is the local triangle-scaling property

$$D_h(x_{k+1}, y_k) \leq \theta_k^{\gamma_k} D_h(z_{k+1}, z_k), \tag{29}$$

where $x_{k+1} = (1 - \theta_k)x_k + \theta_k z_{k+1}$ and $y_k = (1 - \theta_k)x_k + \theta_k z_k$. According to the proof of Lemma 2, we can also use the inequality (21) as stopping criterion, which is implied by (29) and the relatively smooth assumption. For convergence analysis, we only need (21) to hold, which can be less conservative than (29). In Algorithm 2, we use the following inequality as the stopping criterion

$$f(x_{k+1}) \le f(y_k) + \langle \nabla f(y_k), x_{k+1} - y_k \rangle + \theta_k^{\gamma_k} L D_h(z_{k+1}, z_k),$$

which is equivalent to (21) (by subtracting $\Psi(x_{k+1})$ from both sides of the inequality). In practice, this condition often leads to much faster convergence than using (29). Computationally, it is slightly more expensive since it needs to evaluate $f(x_{k+1})$ in addition to $\nabla f(y_k)$ during each inner loop, while (29) does not.

The lower bound $\gamma_{\min}$ can be any known uniform TSE, which guarantees that the stopping criterion can always be satisfied. Without such prior information, we can simply set $\gamma_{\min} = 0$. Since $\gamma_{k+1} \le \gamma_k$ and $\theta_{k+1} \in (0, 1)$, we always have $\theta_{k+1}^{\gamma_{k+1}} \ge \theta_{k+1}^{\gamma_k}$. Therefore

$$\frac{1 - \theta_{k+1}}{\theta_{k+1}^{\gamma_{k+1}}} \le \frac{1 - \theta_{k+1}}{\theta_{k+1}^{\gamma_k}} \le \frac{1}{\theta_k^{\gamma_k}}.$$

By replacing inequality (18) with the one above and repeating the analysis in Section 3.1, we obtain the following result.

**Theorem 2.** *Suppose Assumption A holds, $f$ is $L$-smooth relative to $h$ on $C$, and $\gamma_{\min}$ is a uniform TSE of $D_h$. Then the sequences generated by Algorithm 2 satisfy, for any $x \in \mathrm{dom}\, h$,*

$$F(x_{k+1}) - F(x) \le \left( \frac{\gamma_k}{k + \gamma_k} \right)^{\gamma_k} L D_h(x, x_0), \qquad \forall\, k \ge 0.$$

The convergence rate of ABPG-e is determined by the last value $\gamma_k$. Since we only need to satisfy the local triangle-scaling property (29) instead of the uniform condition (10), it is very likely that $\gamma_k$ is greater than the largest uniform TSE. However, according to Proposition 1, when $k \to \infty$, the limit of $\gamma_k$ (which always exists) cannot be larger than the intrinsic TSE $\gamma_{\mathrm{in}} = 2$. In any case, $\gamma_k$ itself is a *numerical certificate* of an empirical convergence rate of $O(k^{-\gamma_k})$, albeit one that depends on $k$. We always have $\gamma_k \ge \max\{\gamma - \delta, \gamma_{\min}\}$ where $\gamma$ is the uniform TSE. In our numerical experiments in Section 6, the numerical certificate $\gamma_k$ is mostly close to $\gamma_{\mathrm{in}} = 2$.

Compared with ABPG, each iteration of ABPG-e may invoke an inner loop that requires additional computation. However, for finishing the same number of $k$ iterations, the number of extra steps performed by ABPG-e is at most $(\gamma_0 - \gamma_{\min})/\delta$. In practice, we always pick $\delta \ge 0.1$, thus the number of extras steps is a constant of at most a few tens, regardless of the number of iterations $k$.

## 4  ABPG methods with gain adaptation

In this section, we present and analyze an adaptive ABPG method based on the concept of intrinsic TSE developed in Section 2.1. Instead of searching for the largest uniform TSE as in Algorithm 2, we can replace line 2 in Algorithm 1 by

$$z_{k+1} = \arg\min_{z \in C} \left\{ \ell(z|y_k) + G_k \theta_k^{\gamma - 1} L D_h(z, z_k) \right\}$$

---

**Algorithm 3:** ABPG method with gain adaptation (ABPG-g)

---

**input:** initial points $x_0 \in C$, $\gamma > 1$, $\rho > 1$ and $G_{\min} > 0$.

**initialize:** $z_0 = x_0$, $\theta_0 = 1$ and $G_{-1} = 1$

**for** $k = 0, 1, 2, \ldots$ **do**

$\quad M_k = \max\{G_{k-1}/\rho, \ G_{\min}\}$

$\quad$ **repeat** for $t = 0, 1, 2, \ldots$

$\quad\quad G_k = M_k \rho^t$

$\quad\quad$ **if** $k > 0$ **then** compute $\theta_k$ by solving $\dfrac{1 - \theta_k}{G_k \theta_k^\gamma} = \dfrac{1}{G_{k-1} \theta_{k-1}^\gamma}$

$\quad\quad y_k = (1 - \theta_k) x_k + \theta_k z_k$

$\quad\quad z_{k+1} = \arg\min_{z \in C} \left\{ \ell(z | y_k) + G_k \theta_k^{\gamma - 1} L D_h(z, z_k) \right\}$

$\quad\quad x_{k+1} = (1 - \theta_k) x_k + \theta_k z_{k+1}$

$\quad$ **until** $f(x_{k+1}) \leq f(y_k) + \langle \nabla f(y_k), x_{k+1} - y_k \rangle + G_k \theta_k^\gamma L D_h(z_{k+1}, z_k)$

**end**

---

and adjust the additional gain $G_k$ while keeping $\gamma = \gamma_{\text{in}}$ fixed.

Algorithm 3 is such a method with gain adaptation. During each iteration, the algorithm uses an inner loop to search for an value of $G_k$ that satisfies

$$f(x_{k+1}) \leq f(y_k) + \langle \nabla f(y_k), x_{k+1} - y_k \rangle + G_k \theta_k^\gamma L D_h(z_{k+1}, z_k), \tag{30}$$

which is true if the following local triangle-scaling property holds:

$$D_h(x_{k+1}, y_k) = D_h\big((1 - \theta)x_k + \theta z_{k+1}, (1 - \theta)x_k + \theta z_k\big) \leq G_k \theta^\gamma D_h(z_{k+1}, z_k). \tag{31}$$

By definition of the intrinsic TSE, such a $G_k$ (as a function of $x_{k+1}$, $z_k$ and $z_{k+1}$) always exists for $\gamma = \gamma_{\text{in}}$, i.e., the stopping criterion for gain adaptation in Algorithm 3 can always be satisfied. In order to obtain fast convergence, we want the value of $G_k$ to be as small as possible while still satisfying (30). Therefore, at the beginning of each iteration $k$, we always try a tentative gain that is no larger than $G_{k-1}$: $M_k = \max\{G_{k-1}/\rho, \ G_{\min}\}$ with $\rho > 1$. The gain adaptation loop finds the smallest integer $t \geq 0$ such that $G_k = M_k \rho^t$ satisfies the inequality (30).

Another major difference between Algorithm 3 and the previous variants of ABPG is that the sequence $\{\theta_k\}_{k \in \mathbb{N}}$ in Algorithm 3 is generated by solving the equation

$$\frac{1 - \theta_{k+1}}{G_{k+1} \theta_{k+1}^\gamma} = \frac{1}{G_k \theta_k^\gamma}. \tag{32}$$

While an obvious choice is to set $G_{\min} = 1$, we usually set it to be much smaller (say $G_{\min} = 10^{-3}$), which allows the algorithm to converge much faster. Since we don't have a priori upper bound on $G_k$, it is hard to characterize how fast $\theta_k$ converges to zero. In fact, $\{\theta_k\}_{k \in \mathbb{N}}$ may not be a monotone decreasing sequence. Instead of tracking $G_k$ and $\theta_k$ separately, we analyze the convergence of the combined quantity $G_k \theta_k^\gamma$. The following simple lemma will be very useful.

**Lemma 5.** *For any $\alpha, \beta > 0$ and $\gamma \geq 1$, the following inequality holds:*

$$\alpha^\gamma - \beta^\gamma \ \leq \ \gamma(\alpha - \beta)\alpha^{\gamma - 1}.$$

*Proof.* The case of $\gamma = 1$ is obvious. Assume $\gamma > 1$. The desired inequality is equivalent to

$$\alpha^{\gamma-1}\beta \leq \frac{(\gamma-1)\alpha^\gamma + \beta^\gamma}{\gamma} = \frac{(\gamma-1)\alpha^\gamma + 1 \cdot \beta^\gamma}{(\gamma-1)+1}.$$

Applying the weighted arithmetic and geometric mean inequality (26), we have

$$\frac{(\gamma-1)\alpha^\gamma + 1 \cdot \beta^\gamma}{(\gamma-1)+1} \geq \left((\alpha^\gamma)^{\gamma-1}(\beta^\gamma)^1\right)^{\frac{1}{\gamma}} = \alpha^{\gamma-1}\beta,$$

which completes the proof. $\qquad\square$

**Theorem 3.** *Suppose Assumption A holds, $f$ is L-smooth relative to $h$ on $C$, and $\gamma = \gamma_{\mathrm{in}}$ is the intrinsic TSE of $D_h$. Then the sequences generated by Algorithm 3 satisfy, for any $x \in \mathrm{dom}\, h$,*

$$F(x_{k+1}) - F(x) \leq \left(\frac{\gamma}{k+\gamma}\right)^\gamma \overline{G}_k LD_h(x, x_0), \qquad \forall k \geq 0, \tag{33}$$

*where $\overline{G}_k$ is a weighted geometric mean of the gains at each step:*

$$\overline{G}_k = (G_0^\gamma G_1 \cdots G_k)^{\frac{1}{k+\gamma}}. \tag{34}$$

*Proof.* We follow the same steps as in Section 3.1. In light of (31), the inequality (21) becomes

$$F(x_{k+1}) \leq \ell(x_{k+1}|y_k) + G_k\theta_k^\gamma LD_h(z_{k+1}, z_k), \tag{35}$$

and the inequality (23) becomes

$$\frac{1}{G_k\theta_k^\gamma}\left(F(x_{k+1}) - F(x)\right) + LD_h(x, z_{k+1}) \leq \frac{1-\theta_k}{G_k\theta_k^\gamma}\left(F(x_k) - F(x)\right) + LD_h(x, z_k). \tag{36}$$

Plugging in the equality (32), we obtain

$$\frac{1-\theta_{k+1}}{G_{k+1}\theta_{k+1}^\gamma}\left(F(x_{k+1}) - F(x)\right) + LD_h(x, z_{k+1}) \leq \frac{1-\theta_k}{G_k\theta_k^\gamma}\left(F(x_k) - F(x)\right) + LD_h(x, z_k). \tag{37}$$

Then the same arguments in the proof of Theorem 1 lead to

$$F(x_{k+1}) - F(x) \leq G_k\theta_k^\gamma LD_h(x, x_0). \tag{38}$$

Next we derive an upper bound for $G_k\theta_k^\gamma$. For convenience, let's define for $k = 0, 1, 2, \ldots,$

$$A_k = \frac{1}{G_k\theta_k^\gamma}, \qquad a_{k+1} = \frac{1}{G_{k+1}\theta_{k+1}^{\gamma-1}}.$$

Then (32) implies $a_{k+1} = A_{k+1} - A_k$. Moreover, we have

$$A_{k+1} = \frac{1}{G_{k+1}\theta_{k+1}^\gamma} = G_{k+1}^{\frac{1}{\gamma-1}} a_{k+1}^{\frac{\gamma}{\gamma-1}} = G_{k+1}^{\frac{1}{\gamma-1}}\left(A_{k+1} - A_k\right)^{\frac{\gamma}{\gamma-1}}. \tag{39}$$

16

Applying Lemma 5 with $\alpha = A_{k+1}^{1/\gamma}$ and $\beta = A_k^{1/\gamma}$, we obtain

$$A_{k+1} - A_k = \left(A_{k+1}^{\frac{1}{\gamma}}\right)^\gamma - \left(A_k^{\frac{1}{\gamma}}\right)^\gamma \leq \gamma\left(A_{k+1}^{\frac{1}{\gamma}} - A_k^{\frac{1}{\gamma}}\right)A_{k+1}^{\frac{\gamma-1}{\gamma}}.$$

Combining with (39) yields

$$A_{k+1} = G_{k+1}^{\frac{1}{\gamma-1}}\left(A_{k+1} - A_k\right)^{\frac{\gamma}{\gamma-1}} \leq G_{k+1}^{\frac{1}{\gamma-1}}\gamma^{\frac{\gamma}{\gamma-1}}\left(A_{k+1}^{\frac{1}{\gamma}} - A_k^{\frac{1}{\gamma}}\right)^{\frac{\gamma}{\gamma-1}}A_{k+1}$$

We can eliminate the common factor $A_{k+1}$ on both sides of the above inequality to obtain

$$1 \leq G_{k+1}^{\frac{1}{\gamma-1}}\gamma^{\frac{\gamma}{\gamma-1}}\left(A_{k+1}^{\frac{1}{\gamma}} - A_k^{\frac{1}{\gamma}}\right)^{\frac{\gamma}{\gamma-1}},$$

which implies

$$A_{k+1}^{\frac{1}{\gamma}} - A_k^{\frac{1}{\gamma}} \geq \frac{1}{\gamma\, G_{k+1}^{1/\gamma}}, \qquad k = 0, 1, 2, \ldots.$$

Summing the above inequality from step 0 to $k-1$ and using $A_0 = 1/G_0$, we have

$$A_k^{\frac{1}{\gamma}} \geq \sum_{t=1}^{k} \frac{1}{\gamma\, G_t^{1/\gamma}} + A_0^{\frac{1}{\gamma}} = \sum_{t=1}^{k} \frac{1}{\gamma\, G_t^{1/\gamma}} + \frac{1}{G_0^{1/\gamma}} = \frac{1}{\gamma}\left(\sum_{t=1}^{k} \frac{1}{G_t^{1/\gamma}} + \frac{\gamma}{G_0^{1/\gamma}}\right).$$

Using the weighted arithmetic and geometric mean inequality (e.g., [18, Section 2.5]) gives

$$\sum_{t=1}^{k} \frac{1}{G_t^{1/\gamma}} + \frac{\gamma}{G_0^{1/\gamma}} \geq (k+\gamma)\left(\left(\frac{1}{G_0^{1/\gamma}}\right)^\gamma \frac{1}{G_1^{1/\gamma}} \cdots \frac{1}{G_k^{1/\gamma}}\right)^{\frac{1}{k+\gamma}} = (k+\gamma)\left(G_0^\gamma G_1 \cdots G_k\right)^{\frac{-1}{\gamma(k+\gamma)}}.$$

Combining the last two inequalities above, we arrive at

$$A_k \geq \left(\frac{k+\gamma}{\gamma}\right)^\gamma (G_0^\gamma G_1 \cdots G_k)^{\frac{-1}{k+\gamma}}.$$

Therefore,

$$G_k\theta_k^\gamma = \frac{1}{A_k} \leq \left(\frac{\gamma}{k+\gamma}\right)^\gamma (G_0^\gamma G_1 \cdots G_k)^{\frac{1}{k+\gamma}}.$$

Finally, substituting the inequality above into (38) gives the desired result. $\qquad\square$

We note that the geometric mean $\overline{G}_k$ in (34) can be much smaller than the average (arithmetic mean) of $\{G_0, G_1, \ldots, G_k\}$. Under the assumption of uniform Lipschitz smoothness (3), Nesterov [27] proposed an accelerated gradient method with non-monotone line search. However, the complexity obtained there still depends on the global Lipschitz constant $L$, more specifically, replacing $\overline{G}_k L$ in (33) with $\rho L$ when $\gamma = 2$. Our result in (33) can be tighter if the local Lipschitz constants are smaller than $L$ (equivalently with $G_k < 1$).

**Total number of oracle calls.** In order to estimate the overhead of the gain-adaptation procedure, we follow the approach of [27, Lemma 4]. Notice that each inner loop needs to call a gradient oracle to compute $\nabla f(y_k)$, and also $f(x_{k+1})$ when we use (30) as the stopping criterion for gain adaptation. Let $n_i \geq 1$ be the number of calls of the oracle (for $\nabla f(y_k)$) at the $i$th iteration, for $i = 0, \ldots, k$. Then

$$G_{i+1} = \max\{G_i/\rho, G_{\min}\}\rho^{n_i-1} \geq G_i\rho^{n_i-2}, \qquad i = 0, \ldots, k-1.$$

Thus

$$n_i \leq 2 + \log_\rho \frac{G_{i+1}}{G_i} = 2 + \frac{1}{\ln \rho} \ln \frac{G_{i+1}}{G_i}.$$

Therefore, the total number of oracle calls is

$$N_k = \sum_{i=0}^{k} n_i \leq \sum_{i=0}^{k} \left(2 + \frac{1}{\ln \rho} \ln \frac{G_{i+1}}{G_i}\right) = 2(k+1) + \frac{1}{\ln \rho} \ln \frac{G_k}{G_0}.$$

Roughly speaking, on average each iteration need two oracle calls (unless $G_k$ becomes very large).

**An explicit update rule for $\theta_k$.** As an alternative to calculating $\theta_{k+1}$ by solving the equation (32), we can also use the following explicit update rule:

$$\frac{1}{\theta_{k+1}} = \frac{\gamma\alpha_k}{1 + \alpha_k(\gamma-1)} \frac{1}{\theta_k} + \frac{1}{1 + \alpha_k(\gamma-1)},$$

where $\alpha_k = G_{k+1}/G_k$ for $k = 0, 1, 2, \ldots$. This recursion is obtained by solving a linearized equation of (32). In particular, if $\alpha_k = 1$ for all $k \geq 0$, then this formula produces $\theta_k = \gamma/(k+\gamma)$. The sequence $\{\theta_k\}_{k\in\mathbb{N}}$ generated this way satisfies an inequality obtained by replacing the "=" sign with "$\leq$" in (32). Although Theorem 3 does not apply to this sequence, it often has comparable or even faster performance in practice, especially when the $\alpha_k$'s are close to 1.

## 4.1 Towards the $O(k^{-2})$ convergence rate

Proposition 1 shows that the intrinsic TSE $\gamma_{\rm in} = 2$ for all Bregman distances $D_h$ where $h$ is convex and twice continuously differentiable. This covers most Bregman distances of practical interest. If we run the ABPG-g method (Algorithm 3) with $\gamma = 2$, then Theorem 3 states that the convergence rate is $O(\overline{G}_k k^{-2})$. In order to obtain the $O(k^{-2})$ convergence rate, we need $\overline{G}_k$ to be $O(1)$. In this subsection, we discuss the asymptotic behavior of $G_k$, which dominate the geometric mean $\overline{G}_k$ when $k \to \infty$.

We focus on the concrete case of KL divergence, which is a widely used in first-order optimization algorithms. For the KL divergence, $h(x) = \sum_{i=1}^{n} x^{(i)} \log x^{(i)}$ and $\nabla^2 h(x) = \operatorname{diag}\left(\frac{1}{x^{(1)}}, \ldots, \frac{1}{x^{(n)}}\right)$, thus according to the derivations in Section 2.2,

$$G_{\theta_k}(x_k, z_k, z_{k+1}) \leq \lambda_{\max}\left(\nabla^2 h(v_k)^{-1/2} \nabla^2 h(u_k) \nabla^2 h(v_k)^{-1/2}\right) = \max_{i\in\{1,\ldots,n\}} \frac{v_k^{(i)}}{u_k^{(i)}}, \tag{40}$$

where $u_k \in \left[(1-\theta_k)x_k + \theta_k z_k, (1-\theta_k)x_k + \theta_k z_{k+1}\right]$ and $v_k \in \left[z_k, z_{k+1}\right]$. Suppose the sequence $\{x_k\}$ converges to the optimal solution $x_\star$ and $\theta_k \to 0$, then have $u_k \to x_\star$. If $x_\star$ is an interior point of

the positive orthant or the simplex, meaning $x_\star^{(i)} > 0$ for all coordinates $i$, Then we see from (40) that the bound on $G_{\theta_k}(x_k, z_k, z_{k+1})$ depends on how close $x_\star$ is close to the boundary (assuming $v_k$ is bounded).

The most interesting case is when the optimal solution $x_\star$ is on the boundary, i.e., when $x_\star^{(i)} = 0$ for some coordinates $i$. In fact, in our numerical examples on the D-optimal design problem and Poisson linear inverse problem, most of the solutions are on the boundary. However, we emphasize that the iterates generated by the ABPG algorithm is never on the boundary, but may only converge to the boundary; see Assumption A, especially A.5. Our analysis applies to this case as well. In particular, we can show that if both sequences $\{x_k\}$ and $\{z_k\}$ converge, then $x_k^{(i)} \to 0$ implies $z_k^{(i)} \to 0$ and the convergence rate of $x_k^{(i)}$ is no faster than that of $z_k^{(i)}$. (Here $x_k^{(i)} \to 0$ means $\lim_{k\to\infty} x_k^{(i)} = 0$.) More precisely, we have the following lemma.

**Lemma 6.** *Suppose an algorithm generates two sequences $\{x_k\}$ and $\{z_k\}$ in the strictly positive orthant, satisfying $x_0 = z_0$ and $x_{k+1} = (1 - \theta_k)x_k + \theta_k z_{k+1}$ for all $k \geq 0$. Then*

  *(a) If $\{x_k\}$ converges and $x_k^{(i)} \to 0$ for some coordinate $i$, then there must exists an subsequence of $\{z_k^{(i)}\}$ that converges to 0.*

  *(b) Suppose both sequences $\{x_k\}$ and $\{z_k\}$ converge. If $x_k^{(i)} \to 0$ for some $i$, then it converges at a rate that is no faster than $z_k^{(i)}$ in the following sense: For any monotone decreasing sequence $\{r_k\}$ that converges to 0 and satisfies $z_k^{(i)} \geq r_k$ for all $k \geq 0$, we have $x_k^{(i)} \geq r_k$ for all $k \geq 0$. In particular, we can choose $r_k$ to be the monotone lower envelop of $z_k^{(i)}$, i.e., $r_k = \min\{z_0^{(i)}, z_1^{(i)}, \ldots, z_k^{(i)}\}$.*

*Proof.* By the update rule $x_{k+1} = (1 - \theta_k) + \theta_k z_{k+1}$, we know that each $x_k$ is a convex combination of the points $\{z_0 = x_0, z_1, \ldots, z_k\}$, which all lie in the strictly positive orthant. Therefore,

$$x_k^{(i)} \geq \min\{z_0^{(i)}, z_1^{(i)}, \ldots, z_k^{(i)}\} > 0, \qquad \forall i, k.$$

*Part (a).* Suppose $x_k^{(i)} \to 0$ but there is no subsequence of $\{z_k^{(i)}\}$ converging to 0. Then there must exist an $\epsilon > 0$ such that $z_k^{(i)} > \epsilon$ for all $k \geq 0$. Since $x_k$ is a convex combination of $\{z_0, z_1, \ldots, z_k\}$, this implies

$$x_k^{(i)} \geq \min\{z_0^{(i)}, z_1^{(i)}, \ldots, z_k^{(i)}\} > \epsilon, \qquad \forall k \geq 0,$$

which contradicts with the assumption that $x_k^{(i)} \to 0$. Therefore, there must exists an subsequence of $\{z_k^{(i)}\}$ that converges to 0.

*Part (b).* Suppose $\{r_k\}$ is monotone decreasing and converges to 0. If $z_k^{(i)} \geq r_k$ for all $k \geq 0$, then

$$x_k^{(i)} \geq \min\{z_0^{(i)}, z_1^{(i)}, \ldots, z_k^{(i)}\} \geq \min\{r_0, r_1, \ldots, r_k\} = r_k, \qquad \forall k \geq 0.$$

In this sense, $x_k^{(i)} \to 0$ at a rate that is no faster than $z_k^{(i)}$. $\qquad\qquad\square$

According to the bound in (40) and Lemma 6, if both sequences $\{x_k\}$ and $\{z_k\}$ converge, then

$$G_{\theta_k}(x_k, z_k, z_{k+1}) \leq \max_{i\in\{1,\ldots,n\}} \frac{v_k^{(i)}}{u_k^{(i)}} \approx \max_{i\in\{1,\ldots,n\}} \frac{z_k^{(i)}}{x_k^{(i)}},$$

which can be bounded by a constant, since $x_k^{(i)} \to 0$ at a rate that is no faster than $z_k^{(i)}$. In this case, we have $G_k$ in Algorithm 3 bounded by a constant asymptotically, thus the convergence rate is $O(k^2)$.

For the IS divergence, $h(x) = \sum_{i=1}^n -\log x^{(i)}$ and $\nabla^2 h(x) = \text{diag}\left(\left(\frac{1}{x^{(1)}}\right)^2, \ldots, \left(\frac{1}{x^{(n)}}\right)^2\right)$, thus the situation is very similar to the KL divergence.

However, we are not able to prove the convergence of the sequences $\{x_k\}$ and $\{z_k\}$ without additional assumptions (such as relative strong convexity). Indeed, to the best of our knowledge, convergence of these sequences have not been established even under the classical uniform Lipschitz condition. Therefore, an a priori theoretical guarantee of the $O(k^{-2})$ rate seems to be out of reach in general, which seems to coroborate the recent result in [15] that the $O(k^{-1})$ rate cannot be improved in in general for the class of relatively smooth functions.

Nevertheless, we would like to reiterate the remarks at the end of Sections 1.1. In particular, the class of relatively smooth functions is very large, and the lower bound in [15] is established with a worst-case function with pathological nonsmooth behavior. In practical applications, we always work with one particular reference function which may possess structural properties that allow fast convergence. In Algorithm 3, the sequence $\{G_k\}$ is readily available as part of the computation and we can easily check the magnitude of $\overline{G}_k$. Whenever it is small, we obtain a numerical certificate that the algorithm did converge with the $O(k^{-2})$ rate. This is exactly what we observe in the numerical experiments in Section 6.

# 5    Accelerated Bregman dual averaging method

In this section, we present an accelerated Bregman dual averaging (ABDA) method under the relative smoothness assumption. This method extends Nesterov's accelerated dual averaging method ([26] and [33, Algorithm 3]) to the relatively smooth setting. Here we focus on a simple variant in Algorithm 4 based on the uniform triangle-scaling property, although it is also possible to develop more sophisticated variants with automatic exponent or gain adaptation.

In Algorithm 4, Line 2 defines a sequence of functions $\{\psi_k\}_{k \in \mathbb{N}}$ starting with $\psi_0 \equiv 0$:

$$\psi_{k+1}(x) := \psi_k(x) + \theta_k^{1-\gamma} \ell(x|y_k). \tag{41}$$

In other words, $\psi_{k+1}$ is a weighted sum of the lower approximations in (17) constructed at $y_0, \ldots, y_k$:

$$\psi_{k+1}(x) = \sum_{t=0}^{k} \theta_t^{1-\gamma} \ell(x|y_t). \tag{42}$$

Line 3 in Algorithm 4 can be written as

$$z_{k+1} = \arg\min_{z \in C} \left\{ \langle g_k, z \rangle + \vartheta_k \Psi(z) + Lh(z) \right\} \tag{43}$$

where

$$g_k = \sum_{t=1}^{k} \theta_t^{1-\gamma} \nabla f(y_t), \qquad \vartheta_k = \sum_{t=1}^{k} \theta_t^{1-\gamma}.$$

When implementing Algorithm 4, we only need to keep track of $g_k$ and $\vartheta_k$, and there is no need to maintain the abstract form of $\psi_k(x)$. Here our assumption of $C$ and $\Psi$ being simple means that

20

---

**Algorithm 4:** Accelerated Bregman dual averaging (ABDA) method

---

**input:** initial point $z_0 \in \operatorname{rint} C$ and $\gamma > 1$.

initialize: $x_0 = z_0$, $\psi_0(x) \equiv 0$, and $\theta_0 = 1$.

**for** $k = 0, 1, 2, \ldots$ **do**

  1     $y_k := (1 - \theta_k)x_k + \theta_k z_k$

  2     $\psi_{k+1}(x) := \psi_k(x) + \theta_k^{1-\gamma}\ell(x|y_k)$

  3     $z_{k+1} := \arg\min_{z \in C}\{\psi_{k+1}(z) + Lh(z)\}$

  4     $x_{k+1} := (1 - \theta_k)x_k + \theta_k z_{k+1}$

  5     find $\theta_{k+1} \in (0, 1]$ such that $\frac{1-\theta_{k+1}}{\theta_{k+1}^\gamma} = \frac{1}{\theta_k^\gamma}$

**end**

---

the minimization problem in (43) can be solved efficiently. This requirement is equivalent to that for the BPG method (8) and all variants of the ABPG methods in this paper.

Algorithm 4 (line 5) requires the sequence $\{\theta_k\}_{k \in \mathbb{N}}$ satisfy

$$\frac{1 - \theta_{k+1}}{\theta_{k+1}^\gamma} = \frac{1}{\theta_k^\gamma}, \qquad \forall\, k \geq 0. \tag{44}$$

Under this condition, we can show

$$\vartheta_k = \sum_{i=0}^{k} \theta_i^{1-\gamma} = \frac{1}{\theta_k^\gamma}. \tag{45}$$

To see this, we use induction. Clearly it holds for $k = 0$ if we choose $\theta_0 = 1$. Suppose it holds for some $k \geq 0$, then in light of (45) and (44),

$$\vartheta_{k+1} = \sum_{i=0}^{k+1} \frac{1}{\theta_i^{\gamma-1}} = \frac{1}{\theta_k^\gamma} + \frac{1}{\theta_{k+1}^{\gamma-1}} = \frac{1 - \theta_{k+1}}{\theta_{k+1}^\gamma} + \frac{1}{\theta_{k+1}^{\gamma-1}} = \frac{1 - \theta_{k+1} + \theta_{k+1}}{\theta_{k+1}^\gamma} = \frac{1}{\theta_{k+1}^\gamma}.$$

Therefore the inequality (45) holds for all $k \geq 0$.

To analyze the convergence of Algorithm 4, we need the following simple variant of Lemma 1.

**Lemma 7.** *Suppose $h$ is convex and differentiable on* $\operatorname{rint} C$. *For any closed convex function $\varphi$, if*

$$z = \arg\min_{x \in C}\big\{\varphi(x) + h(x)\big\}$$

*and $h$ is differentiable at $z$, then*

$$\varphi(x) + h(x) \geq \varphi(z) + h(z) + D_h(x, z), \quad \forall\, x \in \operatorname{dom} h.$$

**Lemma 8.** *Suppose Assumption A holds, $f$ is $L$-smooth relative to $h$ on $C$, and $\gamma$ is a uniform TSE of $D_h$. Then the sequences generated by Algorithm 4 satisfy, for all $x \in \operatorname{dom} h$ and all $k \geq 1$,*

$$\frac{1 - \theta_{k+1}}{\theta_{k+1}^\gamma} F(x_{k+1}) - \psi_{k+1}(z_{k+1}) - Lh(z_{k+1}) \;\leq\; \frac{1 - \theta_k}{\theta_k^\gamma} F(x_k) - \psi_k(z_k) - Lh(z_k). \tag{46}$$

21

*Proof.* We can start with the inequality (22):

$$
\begin{aligned}
F(x_{k+1}) &\leq (1-\theta_k)\ell(x_k|y_k) + \theta_k\ell(z_{k+1}|y_k) + \theta_k^\gamma LD_h(z_{k+1},z_k) \\
&= (1-\theta_k)\ell(x_k|y_k) + \theta_k^\gamma\left(\theta_k^{1-\gamma}\ell(z_{k+1}|y_k) + LD_h(z_{k+1},z_k)\right) \\
&\leq (1-\theta_k)F(x_k) + \theta_k^\gamma\left(\theta_k^{1-\gamma}\ell(z_{k+1}|y_k) + LD_h(z_{k+1},z_k)\right).
\end{aligned} \tag{47}
$$

Notice that for $k \geq 1$, $z_k$ is the minimizer of $\psi_k(z) + Lh(z)$ over $C$. We use Lemma 7 to obtain

$$
\psi_k(z_k) + Lh(z_k) + LD_h(z_{k+1},z_k) \leq \psi_k(z_{k+1}) + Lh(z_{k+1}),
$$

which gives

$$
LD_h(z_{k+1},z_k) \leq \psi_k(z_{k+1}) + Lh(z_{k+1}) - \psi_k(z_k) - Lh(z_k). \tag{48}
$$

Combining the inequalities (47) and (48), we obtain

$$
\begin{aligned}
F(x_{k+1}) &\leq (1-\theta_k)F(x_k) + \theta_k^\gamma\left(\theta_k^{1-\gamma}\ell(z_{k+1}|y_k) + \psi_k(z_{k+1}) + Lh(z_{k+1}) - \psi_k(z_k) - Lh(z_k)\right) \\
&= (1-\theta_k)F(x_k) + \theta_k^\gamma\left(\psi_{k+1}(z_{k+1}) + Lh(z_{k+1}) - \psi_k(z_k) - Lh(z_k)\right),
\end{aligned}
$$

where in the last equality we used recursive definition of $\psi_{k+1}$ in (41). Dividing both sides of the above inequality by $\theta_k^\gamma$, we have

$$
\frac{1}{\theta_k^\gamma}F(x_{k+1}) \leq \frac{1-\theta_k}{\theta_k^\gamma}F(x_k) + \psi_{k+1}(z_{k+1}) + Lh(z_{k+1}) - \psi_k(z_k) - Lh(z_k).
$$

Using (44) and rearranging terms gives the desired result (46), which holds for $k \geq 1$. $\qquad\square$

**Theorem 4.** *Suppose Assumption A holds, $f$ is $L$-smooth relative to $h$ on $C$, and $\gamma$ is a uniform TSE of $D_h$. The sequences generated by Algorithm 4 satisfy:*

*(a) If $z_0 = \arg\min_{z\in C} h(z)$, then for any $x \in \operatorname{dom} h$,*

$$
F(x_{k+1}) - F(x) \leq \left(\frac{\gamma}{k+\gamma}\right)^\gamma L\big(h(x) - h(z_0)\big), \qquad \forall k \geq 0; \tag{49}
$$

*(b) Otherwise, for any $x \in \operatorname{dom} h$,*

$$
F(x_{k+1}) - F(x) \leq \left(\frac{\gamma}{k+\gamma}\right)^\gamma L\big(h(x) - h(z_1) + D_h(z_1,z_0)\big), \qquad \forall k \geq 0. \tag{50}
$$

*Proof.* If $z_0 = \arg\min_{z\in C} h(z)$, we use the definition $\psi_0 \equiv 0$ to conclude that

$$
z_0 = \arg\min_{z\in C}\big\{\psi_0(z) + Lh(z)\big\}.
$$

In this case, we can extend the result of Lemma 8 to hold for all $k \geq 0$. Applying the inequality (46) for iterations $0, 1, \ldots, k$, we obtain

$$
\frac{1-\theta_{k+1}}{\theta_{k+1}^\gamma}F(x_{k+1}) - \psi_{k+1}(z_{k+1}) - Lh(z_{k+1}) \leq \frac{1-\theta_0}{\theta_0^\gamma}F(x_0) - \psi_0(z_0) - Lh(z_0) = -Lh(z_0),
$$

where we used $\theta_0 = 1$ and $\psi_0 \equiv 0$. Next using (44) and rearranging terms, we have

$$
\begin{aligned}
\frac{1}{\theta_k^\gamma} F(x_{k+1}) &\leq \psi_{k+1}(z_{k+1}) + Lh(z_{k+1}) - Lh(z_0) \\
&\leq \psi_{k+1}(x) + Lh(x) - Lh(x_0) \qquad\qquad (51) \\
&= \sum_{t=0}^{k} \theta_t^{1-\gamma} \ell(x|y_t) + L\big(h(x) - h(z_0)\big) \\
&\leq \sum_{t=0}^{k} \theta_t^{1-\gamma} F(x) + L\big(h(x) - h(z_0)\big) \\
&= \frac{1}{\theta_k^\gamma} F(x) + L\big(h(x) - h(z_0)\big), \qquad\qquad (52)
\end{aligned}
$$

where the second inequality used the fact that $z_{k+1}$ is the minimizer of $\psi_{k+1}(z) + Lh(z)$, the third inequality used $\ell(x|y_t) \leq F(x)$, and the last equality used (45). Rearranging terms of (52) yields

$$
F(x_{k+1}) - F(x) \leq \theta_k^\gamma L\big(h(x) - h(z_0)\big).
$$

According to Lemma 4, we have $\theta_k \leq \frac{\gamma}{k+\gamma}$ if (44) holds, which gives (49).

If $z_0 \neq \arg\min_{z \in C} h(z)$, then we can only apply (46) for $k \geq 1$ to obtain

$$
\begin{aligned}
\frac{1}{\theta_k^\gamma} F(x_{k+1}) - \psi_{k+1}(z_{k+1}) - Lh(z_{k+1}) &\leq \frac{1-\theta_1}{\theta_1^\gamma} F(x_1) - \psi_1(z_1) - Lh(z_1) \\
&= \frac{1}{\theta_0^\gamma} F(x_1) - \theta_0^{1-\gamma} \ell(z_1|y_0) - Lh(z_1) \\
&= F(z_1) - \ell(z_1|z_0) - Lh(z_1) \\
&\leq LD_h(z_1, z_0) - Lh(z_1),
\end{aligned}
$$

where the first equality used (44), the second equality used $\theta_0 = 1$, $y_0 = z_0$ and $x_1 = z_1$, and the last inequality is due to relative smoothness: $F(z_1) \leq \ell(z_1|z_0) + LD_h(z_1, z_0)$. Therefore,

$$
\begin{aligned}
\frac{1}{\theta_k^\gamma} F(x_{k+1}) &\leq \psi_{k+1}(z_{k+1}) + Lh(z_{k+1}) + LD_h(z_1, z_0) - Lh(z_1) \\
&\leq \frac{1}{\theta_k^\gamma} F(x) + L\big(h(x) - h(z_1) + D_h(z_1, z_0)\big),
\end{aligned}
$$

where the last inequality repeats the arguments from (51) to (52). Rearranging terms leads to

$$
F(x_{k+1}) - F(x) \leq \theta_k^\gamma L\big(h(x) - h(z_1) + D_h(z_1, z_0)\big),
$$

and further applying Lemma 4 gives the desired result (50). $\qquad\square$

As a sanity check, we show that the right-hand-side of (50) is strictly positive for any $x \in \operatorname{dom} h$ such that $F(x) < F(z_1) + LD_h(x, z_1)$. We exploit the fact that $z_1 = \arg\min_{z \in C}\{\ell(z|z_0) + Lh(z)\}$. Using Lemma 7, we have

$$
\ell(z_1|z_0) + Lh(z_1) \leq \ell(x|z_0) + Lh(x) - LD_h(x, z_1),
$$

(a) ABPG method with different TSE $\gamma$.



(b) Local triangle scaling gain $\widehat{G}_k$.

Figure 2: D-optimal design: random problem instance with $m = 80$ and $n = 200$.

which implies

$$L\big(h(x) - h(z_1)\big) \geq LD_h(x, z_1) + \ell(z_1|z_0) - \ell(x|z_0).$$

Then we have

$$
\begin{aligned}
L\big(h(x) - h(z_1) + D_h(z_1, z_0)\big) &\geq LD_h(x, z_1) + \ell(z_1|z_0) - \ell(x|z_0) + LD_h(z_1, z_0) \\
&= LD_h(x, z_1) + \big(\ell(z_1|z_0) + LD_h(z_1, z_0)\big) - \ell(x|z_0) \\
&\geq LD_h(x, z_1) + F(z_1) - \ell(x|z_0) \\
&\geq LD_h(x, z_1) + F(z_1) - F(x),
\end{aligned}
$$

where the second inequality used the upper bound in (17), and the last inequality used the lower bound in (17). Therefore, for any $x$ such that $F(x) < F(z_1) + LD_h(x, z_1)$, we have

$$L\big(h(x) - h(z_1) + D_h(z_1, z_0)\big) > LD_h(x, z_1) \geq 0.$$

This completes the proof.

# 6 Numerical experiments

We consider three applications of relatively smooth convex optimization: D-optimal experiment design, Poisson linear inverse problem, and relative-entropy nonnegative regression. For each application, we compare the algorithms developed in this paper with the BPG method (8) and demonstrate significant performance improvement. Our implementations and experiments are shared through an open-source repository at `https://github.com/linxiaolx/accbpg`.

## 6.1 D-optimal experiment design

Given $n$ vectors $v_1, \ldots, v_n \in \mathbb{R}^m$ where $n \geq m + 1$, the D-optimal design problem is

$$
\begin{aligned}
\text{minimize} \quad & f(x) := -\log \det \big(\textstyle\sum_{i=1}^{n} x^{(i)} v_i v_i^T\big) \\
\text{subject to} \quad & \textstyle\sum_{i=1}^{n} x^{(i)} = 1 \\
& x^{(i)} \geq 0, \quad i = 1, \ldots, n.
\end{aligned}
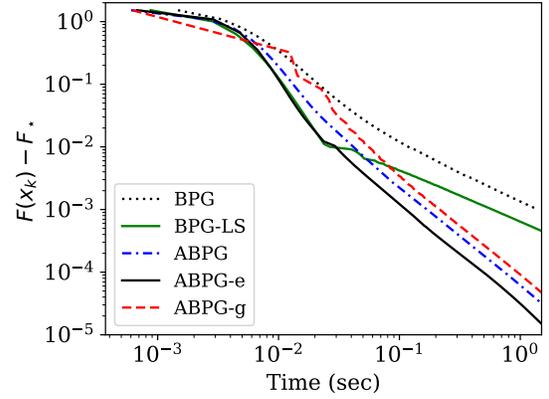\tag{53}
$$

24

(a) Objective gap versus iterations (linear scale).



(b) Objective gap versus iterations (log scale).



(c) Objective gap versus CPU time (linear scale).



(d) Objective gap versus CPU time (log scale).

Figure 3: D-optimal design: random instance with $m = 80$ and $n = 200$; $\gamma = 2$ for all ABPG variants.

In the form of problem (1), we have $\Psi \equiv 0$, $F(x) \equiv f(x)$, and $C$ is the standard simplex in $\mathbb{R}^n$. In statistics, this problem corresponds to maximizing the determinant of the Fisher information matrix (e.g., [19, 1]). It is shown in [21] that $f$ defined in (53) is 1-smooth relative to Burg's entropy $h(x) = -\sum_{i=1}^{n} \log(x^{(i)})$ on $\mathbb{R}_+^n$. In this case, $D_h$ is the IS-distance defined in (12).

### 6.1.1 Experiment on synthetic data

In our first experiment, we set $m = 80$ and $n = 200$ and generated $n$ random vectors in $\mathbb{R}^m$, where the entries of the vectors were generated following independent Gaussian distributions with zero mean and unit variance. The results are shown in Figures 2 and 3.

Figure 2(a) shows the reduction of optimality gap by the BPG method (8) and the ABPG method (Algorithm 1) with four different values of $\gamma$. For $\gamma = 1$, the ABPG method converges with $O(k^{-1})$ rate, but is slower than the BPG method. When we increase $\gamma$ to 1.5 and then 2, the ABPG method is significantly faster than BPG. Interestingly, ABPG still converges with $\gamma = 2.2$ (which is larger than the intrinsic TSE $\gamma_{\text{in}} = 2$) and is even faster than with $\gamma = 2$. To better
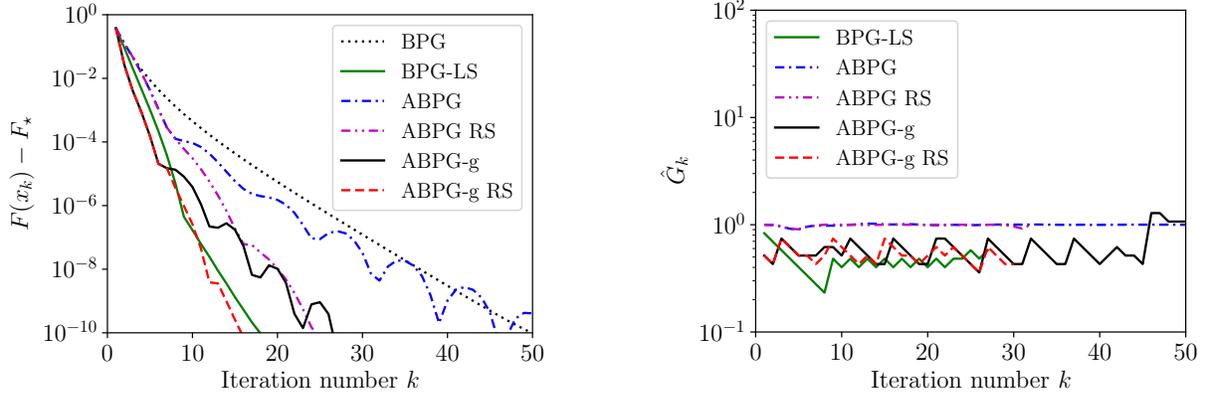
25

Figure 4: D-optimal design: random problem instance with $m = 80$ and $n = 120$.

understand this phenomenon, we plot the local triangle-scaling gain

$$\widehat{G}_k = \frac{D_h(x_{k+1}, y_k)}{\theta^\gamma D_h(z_{k+1}, z_k)} = \frac{D_h((1-\theta)x_k + \theta z_{k+1}, (1-\theta)x_k + \theta z_k)}{\theta^\gamma D_h(z_{k+1}, z_k)}. \tag{54}$$

Figure 2(b) shows that for $\gamma = 1.0$ and $1.5$, $\widehat{G}_k$ is mostly much smaller than 1. For $\gamma = 2$, $\widehat{G}_k$ is much closer to 1 but always less than 1. This gives a numerical certificate that the ABPG method converged with $O(k^{-2})$ rate. For $\gamma = 2.2$, $\widehat{G}_k$ stayed close to 1 for the first 700 iterations and then jumped to 3 and stayed around. The method diverges with larger value of $\gamma$. We didn't plot the ABDA method (Algorithm 4) because it overlaps with ABPG for the same value of $\gamma$ when the initial point is taken as the center of the simplex, see part (a) of Theorem 4.

Figure 3(a) compares the basic BPG and ABPG methods with their adaptive variants. The BPG-LS method is a variant of BPG equipped with the same adaptive line-search scheme in Algorithm 3 (see also [27, Method 3.3]). For all variants of ABPG, we set $\gamma = \gamma_{\text{in}} = 2$. For BPG-LS and ABPG-g, we set $\rho = 1.5$ for adjusting the gain $G_k$. The adaptive variants converged in fewer iterations than their respective basic versions. Figure 3(b) shows the same results in log-log scale. We can clearly see the different slopes of the BPG variants and ABPG variants, demonstrating their $O(k^{-1})$ and $O(k^{-2})$ convergence rates respectively. For ABPG-e, we started with $\gamma_0 = 3$ and it eventually settled down to $\gamma = 2$, which is reflected in its gradual change of slope in Figure 3(b).

We also show the comparison in terms of CPU time in Figures 3(c) and 3(d). As remarked at the end of Section 3.2, the ABPG-e method only take a constant number more iterations than ABPG, thus its their comparison is very similar to the case with number of iterations. For ABPG-g, the analysis in Section 4 on page 18 shows that the number of gradient calls and proximal computations is roughly twice of the ABPG method with the same number of iterations. This is exactly what we observe in Figures 3(c) and 3(d). Given such predictable scaling between number of iterations and CPU time, we only show comparisons in the number of iterations in the rest numerical experiments.

Figure 4 shows the comparison of different methods on another random problem instance with $m = 80$ and $n = 120$. All methods converge much faster and reach very high precision. In particular, BPG and BPG-LS look to have linear convergence. This indicates that this problem instance is much better conditioned and the objective function may be strongly convex relative to Burg's entropy. In this case, it is shown in [21] that the BPG method attains linear convergence. The ABPG and ABPG-g methods demonstrate periodic non-monotone behavior. A well-known

26

technique to avoid such oscillations and attain fast linear convergence is to restart the algorithm whenever the function value starts to increase [30]. We applied restart (RS) to both ABPG and ABPG-g, which resulted in a much faster convergence as shown in Figure 4.

### 6.1.2  Experiment on real data

In our second experiment, we construct D-optimal design instances from LibSVM data [12]. In particular, we consider several regression datasets – the goal is to find the most relevant data points where one shall run the experiment to evaluate the corresponding label.

Figure 5 shows the results on four different datasets: `abalone` ($n = 4177, m = 8$), `bodyfat` ($n = 252, m = 14$), `mpg` ($n = 392, m = 7$) and `housing` ($n = 506, m = 13$). The left column indicates that in each case, the best performance of ABPG is achieved with large TSE $\gamma = 2$ and $\gamma = 2.2$. Furthermore, ABPG with $\gamma > 1$ always compared favorably over plain BPG.

Next, the second column of Figure 5 shows that both ABPG-g and ABPG (with $\gamma = 2$) always significantly outperform BPG and BPG-LS. We have chosen log-log scale of the plot to contrast the $O(k^{-1})$ convergence rate of BPG (with line search) with the $O(k^{-2})$ convergence rate of ABPG and ABPG-g. In the third column, we plot the local triangle-scaling gains. They serve as numerical certificates of the empirical $O(k^{-2})$ convergence rate of ABPG and its variants. In particular, we see that $G_k$ for the ABPG-g algorithm is mostly flat and less than one.

## 6.2  Poisson linear inverse problem

In Poisson inverse problems (e.g., [14, 7]), we are given a nonnegative observation matrix $A \in \mathbb{R}_+^{m \times n}$ and a noisy measurement vector $b \in \mathbb{R}_{++}^m$, and the goal is to reconstruct the signal $x \in \mathbb{R}_+^n$ such that $Ax \approx b$. A natural measure of closeness of two nonnegative vectors is the KL-divergence defined in (11). In particular, minimizing $D_{\mathrm{KL}}(b, Ax)$ corresponds to maximizing the Poisson log-likelihood function. We consider problems of the form

$$\underset{x \in \mathbb{R}_+^n}{\text{minimize}} \ \ F(x) := D_{\mathrm{KL}}(b, Ax) + \Psi(x),$$

where $\Psi(x)$ is a simple regularization function. It is shown in [3] that the function $f(x) = D_{\mathrm{KL}}(b, Ax)$ is $L$-smooth relative to $h(x) = -\sum_{i=1}^n \log(x^{(i)})$ on $\mathbb{R}_+^n$ for any $L \geq \|b\|_1 = \sum_{i=1}^m b^{(i)}$. Therefore, in the BPG and ABPG methods, we use again the IS-distance $D_{\mathrm{IS}}$ defined in (12) as the proximity measure.

Figure 6 shows our computational results for a randomly generated instance with $m = 200$ and $n = 100$ and $\Psi \equiv 0$ (no regularization). The entries of $A$ and $b$ are generated following independent uniform distribution over the interval $[0, 1]$.

Figure 6(a) shows the reduction of objective gap by BPG and ABPG with $\gamma = 1.0$, 1.5 and 2.0, as well as the ABDA method (Algorithm 4). ABPG and ABDA with $\gamma = 2$ mostly overlap each other in this figure. Figure 6(b) plots the same results in log-log scale, which reveals that ABPG and ABDA (both with $\gamma = 2$) behave quite differently in the beginning. The ABDA method has a jump of objective value at $k = 1$ because $z_0 \neq \arg\min_{z \in C} h(z)$, and its convergence rate is governed by part (b) of Theorem 4. In fact, for $C = \mathbb{R}_+^n$, Burg's entropy $h(x) = -\sum_{i=1}^n \log(x^{(i)})$ is unbounded below as $\|x\| \to \infty$. In contrast, for the D-optimal design problem in Section 6.1, $C$ is the standard simplex, and if we choose $z_0 = x_0 = (1/n, \ldots, 1/n)$ then $z_0 = \arg\min_{z \in C} h(z)$. In that case, we can show that ABPG and ABDA are equivalent when $\Psi \equiv 0$.
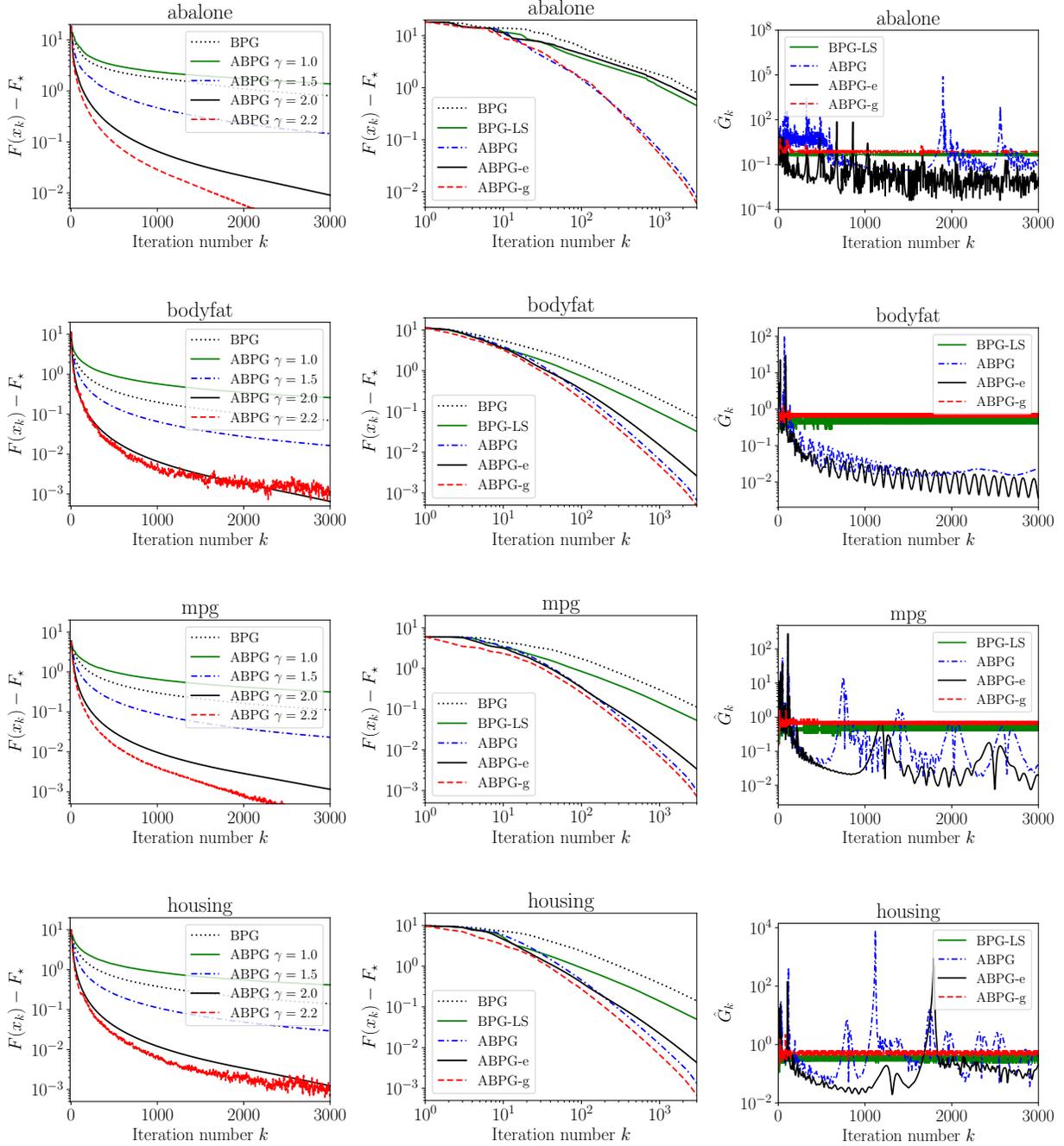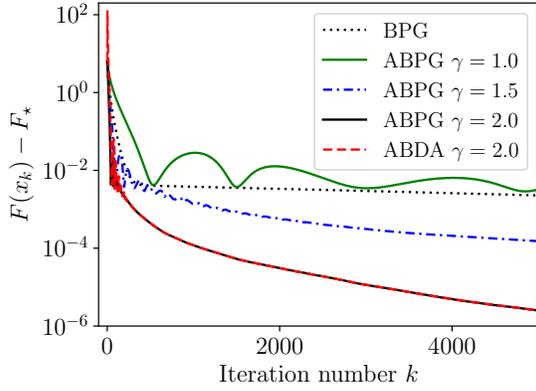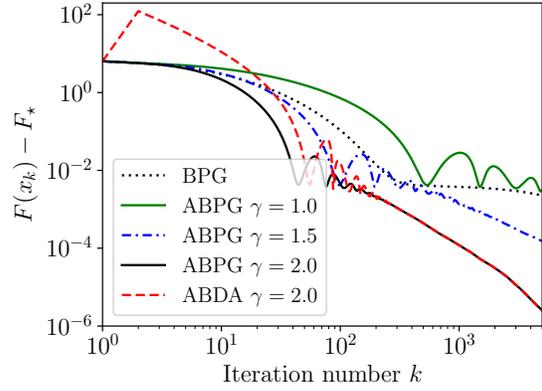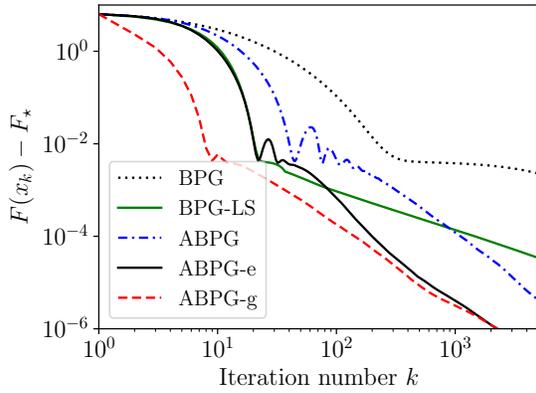
27

Figure 5: D-optimal design on several LibSVM datasets: `abalone`, `bodyfat`, `mpg` and `housing`. The first column compares the BPG method, ABPG method with $\gamma \in \{1.0, 1.5, 2.0, 2.2\}$ and the ABDA method with $\gamma = 2.0$, plotted in log-linear scale. The second column compares BPG, BPG-LS and three ABPG variants with $\gamma = 2$, plotted in log-log scale, and the third column plots the corresponding local triangle-scaling gains. For ABPG and ABPG-e, these are the $\hat{G}_k$ estimated using (54); For BPG-LS and ABPG-g, they are calculated as part of the gain adaptation schemes.
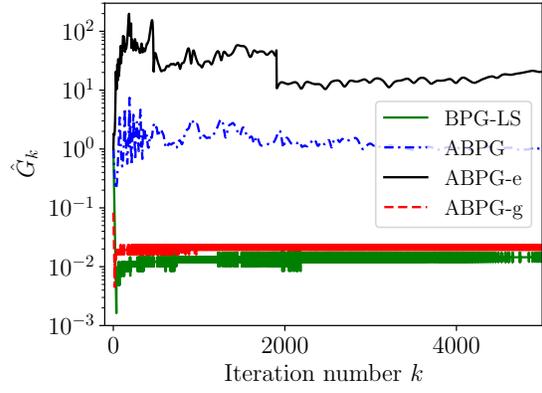
(a) ABPG (varying $\gamma$) and ABDA ($\gamma = 2$) methods.

(b) Same results in (a) in log-log plot.

(c) Adaptive ABPG methods $\gamma = 2$ (log-log plot).

(d) Local triangle-scaling gain $\widehat{G}_k$.

Figure 6: Poisson linear inverse problem: random instance with $m = 200$ and $n = 100$.

Figure 6(c) compares the basic and adaptive variants of BPG and ABPG. For the ABPG and ABPG-g methods, we set $\gamma = \gamma_{\text{in}} = 2$. For ABPG-e, we start with $\gamma_0 = 3$, and the final $\gamma_k = 2.8$ after $k = 5000$ iterations ($\delta = 0.2$ in Algorithm 2). Although ABPG-e uses a much larger $\gamma$ most of the time, we see ABPG-g converges faster than ABPG-e in the beginning and they eventually become similar. This can be explained through the effective triangle-scaling gains plotted in Figure 6(d). For ABPG and ABPG-e, the effective gains plotted are $\widehat{G}_k$ defined in (54). For BPG-LS and ABPG-g, we plot the $G_k$'s which are adjusted directly in the algorithms. For ABPG-g, $G_k \approx 0.025$ most of the time. The effective $\widehat{G}_k$ for ABPG-e is almost 1000 times larger, which counters the large value of $\gamma$ used. The sudden reduction of $\widehat{G}_k$ around $k = 2000$ is when $\gamma$ is reduced from 3 to 2.8. We expect $\gamma_k \to 2$ as $k$ continues to increase.

Figure 7 shows the results for a randomly generated instance with $m = 100$ and $n = 1000$. In this case, since $m < n$, we added a regularization $\Psi(x) = (\lambda/2)\|x\|^2$ with $\lambda = 0.001$. ABPG-g has the best performance. Again we observe that $G_k \ll 1$ most of the time, which gives a numerical certificate that the ABPG methods do converge with $O(k^{-2})$ rate.
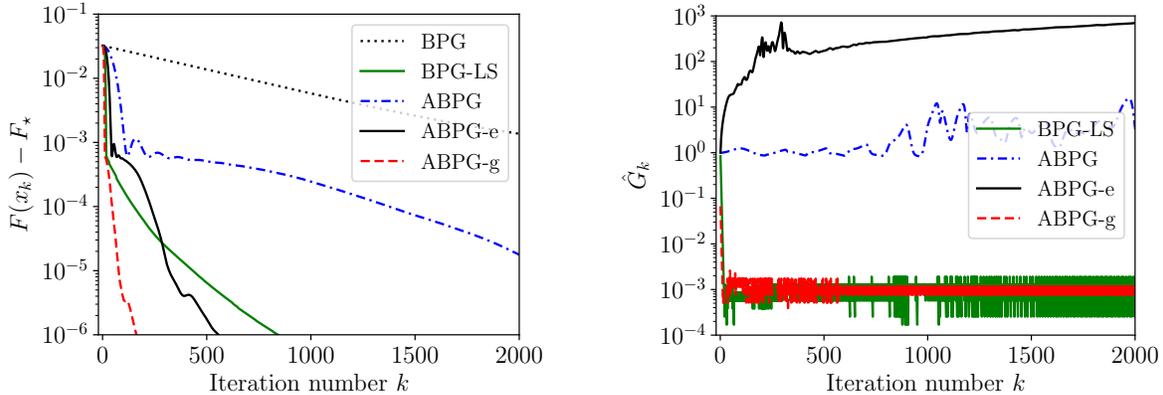
Figure 7: Poisson linear inverse problem: random instance with $m = 100$ and $n = 1000$.

## 6.3 Relative-entropy nonnegative regression

An alternative approach for solving the nonnegative linear inverse problem described in Section 6.2 is to minimize $D_{\mathrm{KL}}(Ax, b)$, i.e.,

$$\underset{x \in \mathbb{R}_+^n}{\operatorname{minimize}} \ F(x) := D_{\mathrm{KL}}(Ax, b) + \Psi(x).$$

In this case, it is shown in [3] that $f(x) = D_{\mathrm{KL}}(Ax, b)$ is $L$-smooth relative to the Boltzmann-Shannon entropy $h(x) = \sum_{i=1}^n x^{(i)} \log(x^{(i)})$ on $\mathbb{R}_+^n$ for any $L$ such that

$$L \ \geq \ \max_{1 \leq j \leq n} \sum_{i=1}^m A_{ij} \ = \ \max_{1 \leq j \leq n} \|A_{\cdot j}\|_1$$

where $A_{\cdot j}$ denotes the $j$th column of $A$. Therefore, in the BPG and ABPG methods, we use the KL-divergence $D_{\mathrm{KL}}$ defined in (11) as the proximity measure. In our experiment, we apply $\ell_1$-regularization $\Psi(x) = \lambda\|x\|_1$ with $\lambda = 0.001$.

Figure 8(a) shows the results for a randomly generated instance with $m = 1000$ and $n = 100$. For all variants of the ABPG method, we set $\gamma = \gamma_{\mathrm{in}} = 2$. Figure 8(b) shows the results for a random instance with $m = 100$ and $n = 1000$. In this case, we clearly see linear convergence of the BPG and BPG-LS methods. Since the accelerated methods demonstrate oscillations in objective value, we tried the restart (RS) trick [30] and obtained faster convergence with apparent linear rate. In contrast, ABPG methods with restart do not make any difference in Figure 8(a). Although not shown here, for the ABPG-g method, we always obtain small gains $G_k \leq 1$ at each step. Therefore their geometric mean $\overline{G}_k$ is also small, which serves as a certificate of the $O(k^{-2})$ convergence rate for this problem instance.

## Acknowledgments

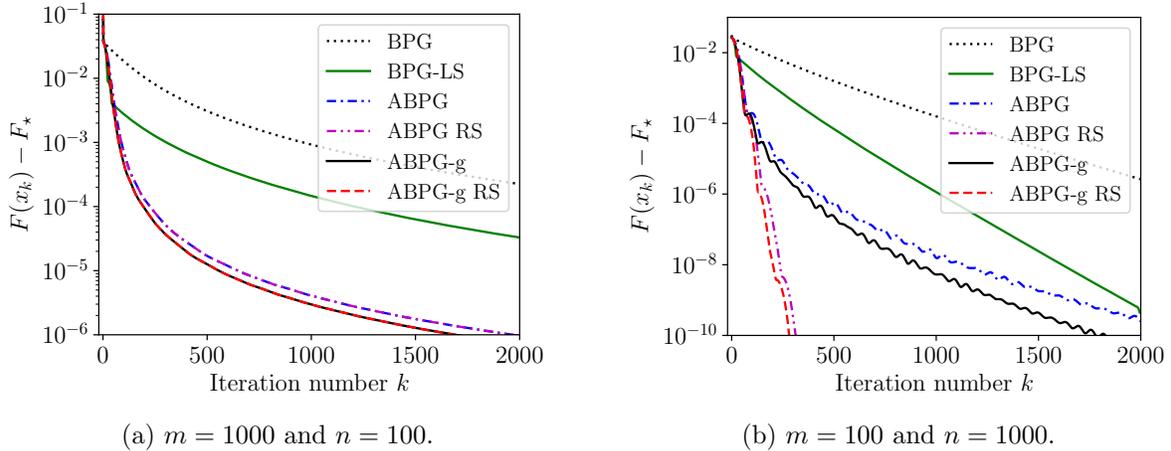(a) $m = 1000$ and $n = 100$.  (b) $m = 100$ and $n = 1000$.

Figure 8: Two random instances of relative entropy nonnegative regression ($\gamma = 2$ for ABPG).

# References

[1] C. L. Atwood. Optimal and efficient designs of experiments. *The Annals of Mathematical Statistics*, 40(5):1570–1602, 1969.

[2] A. Auslender and M. Teboulle. Interior gradient and proximal methods for convex and conic optimization. *SIAM Journal on Optimization*, 16(3):697–725, 2006.

[3] H. H. Bauschke, J. Bolte, and M. Teboulle. A descent Lemma beyond Lipschitz gradient continuity: first-order method revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.

[4] H. H. Bauschke and J. M. Borwein. Joint and separate convexity of the Bregman distance. In D. Butnariu, Y. Censor, and S. Reich, editors, *Inherently Parallel Algorithms in Feasibility and Optimization and their Applications (Haifa 2000)*, pages 23–26. Elsevier, 2001.

[5] A. Beck. *First-Order Methods in Optimization*. MOS-SIAM Series on Optimization. SIAM, 2017.

[6] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.

[7] M. Bertero, P. Boccacci, G. Desiderà, and G. Vicidomini. Image deblurring with Poisson data: from cells to galaxies. *Inverse Problems*, 25(12), 2009.

[8] B. Birnbaum, N. R. Devanur, and L. Xiao. Distributed algorithms via gradient descent for Fisher markets. In *Proceedings of the 12th ACM conference on Electronic Commerce*, pages 127–136, San Jose, California, USA, June 2011.

[9] L. M. Bregman. The relaxation method of finding the common points of convex sets and its application to the solution of problems in convex programming. *U.S.S.R. Comput. Math. and Math. Phys.*, 7:200–217, 1967.

[10] Y. Censor and A. Lent. An iterative row-action method for interval convex programming. *Journal of Optimization theory and Applications*, 34(3):321–353, 1981.

[11] Y. Censor and S. A. Zenios. Proximal minimization algorithm withd-functions. *Journal of Optimization Theory and Applications*, 73(3):451–464, 1992.

[12] C.-C. Chang and C.-J. Lin. Libsvm: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27, 2011.

[13] G. Chen and M. Teboulle. Convergence analysis of a proximal-like minimization algorithm using Bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, August 1993.

[14] I. Csiszár. Why least squares and maximum entropy? an axiomatic approach to inference for linear iverse problems. *The Annals of Statistics*, 19(4):2032–2066, 1991.

[15] R.-A. Dragomir, A. B. Taylor, A. d'Aspremont, and J. Bolte. Optimal complexity and certification of bregman first-order methods. Preprint, arXiv:1911.08510, 2019.

[16] D. H. Gutman and J. F. Peña. Perturbed Fenchel duality and first-order methods. Preprint, arXiv:1812.10198, 2018.

[17] F. Hanzely and P. Richtárik. Fastest rates for stochastic mirror descent methods. arXiv preprint arXiv:1803.07374, 2018.

[18] G. Hardy, J. E. Littlewood, and G. Pólya. *Inequalities*. Cambridge University Press, 2nd edition, 1952.

[19] J. Kiefer and J. Wolfowitz. Optimal design in regression problems. *The Annals of Mathematical Statistics*, 30(2):271–294, 1959.

[20] H. Lu. "Relative-continuity" for non-Lipschitz non-smooth convex optimization using stochastic (or deterministic) mirror descent. *INFORMS Journal on Optimization*, 1(4):288–303, 2019.

[21] H. Lu, R. M. Freund, and Y. Nesterov. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.

[22] A. Nemirovsky and D. Yudin. *Problem Complexity and Method Efficiency in Optimization*. J. Wiley & Sons, New York, 1983.

[23] Y. Nesterov. A method for solving a convex programming problem with convergence rate $O(1/k^2)$. *Soviet Mathematics - Doklady*, 27(2):372–376, 1983.

[24] Y. Nesterov. On an approach to the construction of optimal methods of minimization of smooth convex functions. *Èkonom. i. Mat. Metody*, 24:509–517, 1988.

[25] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer, Boston, 2004.

[26] Y. Nesterov. Smooth minimization of nonsmooth functions. *Mathematical Programming*, 103:127–152, 2005.

[27] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming, Ser. B*, 140:125–161, 2013.

[28] Y. Nesterov. Universal gradient methods for convex optimization problems. *Mathematical Programming, Ser. A*, 152:381–404, 2015.

[29] Y. Nesterov. Implementable tensor methods in unconstrained convex optimization. *Mathematical Programming*, 186:157–183, 2021.

[30] B. O'Donoghue and E. Candès. Adaptive restart for accelerated gradient schemes. *Foundations of Computational Mathematics*, 15(3):715–732, 2015.

[31] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.

[32] M. Teboulle. A simplified view of first order methods for optimization. *Mathematical Programming, Ser. B*, 170:67–96, 2018.

[33] P. Tseng. On accelerated proximal gradient methods for convex-concave optimization. Unpublished manuscript, 2008.

[34] Y. Zhou, Y. Liang, and L. Shen. A simple convergence analysis of Bregman proximal gradient algorithm. *Computational Optimization and Applications*, 93:903–912, 2019.