

Nonconvex Stochastic Bregman Proximal Gradient Method with Application to Deep Learning

Kuangyu Ding

Department of Mathematics

National University of Singapore

10 Lower Kent Ridge Road, Singapore 119076

KUANGYUD@U.NUS.EDU

Jingyang Li

Department of Mathematics

National University of Singapore

10 Lower Kent Ridge Road, Singapore 119076

LI.JINGYANG@U.NUS.EDU

Kim-Chuan Toh

Department of Mathematics

Institute of Operations Research and Analytics

National University of Singapore

10 Lower Kent Ridge Road, Singapore 119076

MATTOHKC@NUS.EDU.SG

Editor: Sanjiv Kumar

Abstract

Stochastic gradient methods for minimizing nonconvex composite objective functions typically rely on the Lipschitz smoothness of the differentiable part, but this assumption fails in many important problem classes like quadratic inverse problems and neural network training, leading to instability of the algorithms in both theory and practice. To address this, we propose a family of stochastic Bregman proximal gradient (SBPG) methods that only require smooth adaptivity. SBPG replaces the quadratic approximation in SGD with a Bregman proximity measure, offering a better approximation model that handles non-Lipschitz gradients in nonconvex objectives. We establish the convergence properties of vanilla SBPG and show it achieves optimal sample complexity in the nonconvex setting. Experimental results on quadratic inverse problems demonstrate SBPG's robustness in terms of stepsize selection and sensitivity to the initial point. Furthermore, we introduce a momentum-based variant, MSBPG, which enhances convergence by relaxing the mini-batch size requirement while preserving the optimal oracle complexity. We apply MSBPG to the training of deep neural networks, utilizing a polynomial kernel function to ensure smooth adaptivity of the loss function. Experimental results on benchmark datasets confirm the effectiveness and robustness of MSBPG in training neural networks. Given its negligible additional computational cost compared to SGD in large-scale optimization, MSBPG shows promise as a universal open-source optimizer for future applications.

Keywords: Nonconvex stochastic algorithm, Bregman distance, Smooth adaptivity, Deep neural network, Algorithmic stability

1. Introduction

In this paper, we present and analyze a family of nonconvex stochastic Bregman proximal gradient methods (SBPG) for solving the following generic stochastic minimization problem:

$$\min_{\mathbf{x} \in \bar{C}} \mathbb{E}_{\boldsymbol{\xi}}[f(\mathbf{x}, \boldsymbol{\xi})] + R(\mathbf{x}), \quad (1)$$

where $f(\cdot, \boldsymbol{\xi})$ is a nonconvex differentiable function on \bar{C} , R is a proper lower-semicontinuous convex function, $\boldsymbol{\xi}$ is a random variable, and \bar{C} is the closure of C , which is a nonempty convex open subset of \mathbb{R}^d . We denote $F(\mathbf{x}) := \mathbb{E}_{\boldsymbol{\xi}}[f(\mathbf{x}, \boldsymbol{\xi})]$, and $\Phi(\mathbf{x}) := F(\mathbf{x}) + R(\mathbf{x})$. This type of stochastic minimization problem is common in machine learning and statistics (Hastie et al., 2009; Shapiro et al., 2021; Zhang, 2004), where the optimizer has limited access to the distribution of $\boldsymbol{\xi}$ and can only draw samples from it. In many instances, the smooth part of the objective function $F(\mathbf{x})$ can be formulated as a finite-sum structure $F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$. Although the distribution of $\boldsymbol{\xi}$ is known in such cases, when n is extremely large, calculating the true gradient for the smooth part of the objective function becomes extremely expensive. As a result, stochastic first-order methods, originating from the work of Robbins and Monro (1951), have emerged as the prevailing approach for solving these large-scale optimization problems. In particular, stochastic (proximal) gradient descent and its numerous variants (Duchi et al., 2011; Duchi and Singer, 2009; Gu et al., 2020; Kingma and Ba, 2014; Allen-Zhu, 2018; Wang et al., 2022) have been widely used in large-scale stochastic optimization for machine learning (LeCun et al., 2015; Shapiro et al., 2021; Zhang, 2004). From a modeling perspective, stochastic (proximal) gradient descent can be viewed as minimizing a sequence of upper quadratic approximations of the nonconvex objective $\Phi(\mathbf{x})$:

$$\mathbf{x}^{k+1} = \operatorname{argmin}_{\mathbf{x} \in \bar{C}} \left\{ \underbrace{F(\mathbf{x}^k, \boldsymbol{\Xi}_k) + \langle \tilde{\nabla}_k, \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}^k\|^2 + R(\mathbf{x})}_{F_{\mathbf{x}^k}(\mathbf{x}): \text{model of } F \text{ at } \mathbf{x}^k} \right\}, \quad (2)$$

where $F(\mathbf{x}^k, \boldsymbol{\Xi}_k) := \frac{1}{|\boldsymbol{\Xi}_k|} \sum_{\boldsymbol{\xi} \in \boldsymbol{\Xi}_k} f(\mathbf{x}^k, \boldsymbol{\xi})$, $\boldsymbol{\Xi}_k$ is the set of samples of $\boldsymbol{\xi}$ at the k -th iteration, and $\tilde{\nabla}_k$ is an estimator of the exact gradient $\nabla F(\mathbf{x}^k)$. This modeling perspective is well-established in deterministic optimization and has been used in methods such as Newton method, Gauss-Newton method, bundle method, and trust-region method, as discussed in various sources such as Hiriart-Urruty and Lemaréchal (1993); Nesterov (2003); Lin et al. (2007); Paren et al. (2022).

Despite their widespread use, stochastic gradient methods (2) face several well-known challenges both in theory and practice. One of the key assumptions underlying the analysis of these methods is the Lipschitz continuity of the gradient of the differentiable part, which is critical for ensuring convergence and establishing complexity results. However, this assumption is not always valid. For instance, even a simple function like $F(x) = x^4$ does not admit a globally Lipschitz gradient over \mathbb{R} , illustrating the limitations in analyzing stochastic gradient methods in more general settings. In addition, choosing the appropriate stepsize is another challenge in both the theoretical analysis and the practical usage

of stochastic gradient methods. In practice, the stepsize plays a decisive role—arguably the most important—in determining the convergence behavior of the algorithm. Finding the optimal stepsize can be time-consuming, as engineers often need to conduct numerous experiments to fine-tune it. From a theoretical perspective, choosing a cautious stepsize is necessary to ensure the descent property at each iteration, which is typically proportional to the inverse of the Lipschitz smoothness parameter. As a result, the absence of Lipschitz smoothness can lead to instability in classical stochastic first-order methods, complicating both their theoretical analysis and practical implementation.

To address these issues, classical approaches often resort to either line search or more complicated inner loops, but these methods can negatively impact the efficiency of the algorithm or even become intractable in a stochastic setting. For example, stochastic proximal point algorithm (PPA) models the approximation of $F(\mathbf{x})$ in (2) as $F_{\mathbf{x}^k}(\mathbf{x}) = F(\mathbf{x}, \boldsymbol{\xi}_k) + \frac{1}{2\alpha_k}\|\mathbf{x} - \mathbf{x}^k\|^2$ (Bertsekas, 2011; Bianchi, 2016; Patrascu and Necoara, 2017; Rockafellar, 1976), making the selection of stepsize α_k more robust than the original model (2). However, the application of stochastic PPA is limited due to the difficulty of solving the subproblems, especially for complicated objectives, such as training deep neural networks. In such cases, solving the subproblem is almost as difficult as solving the original problem, rendering the method impractical. Recently, Bauschke et al. (2017); Lu et al. (2018) have proposed using Bregman proximity measures to relax the assumption of gradient Lipschitz continuity to smooth adaptivity. The Bregman gradient method was first introduced as the mirror descent scheme by Nemirovskij and Yudin (1983) for minimizing convex nonsmooth functions. From the modeling perspective, Bregman methods consider the following subproblem at each iteration:

$$\mathbf{x}^{k+1} = \underset{\mathbf{x} \in \bar{C}}{\operatorname{argmin}} \left\{ \underbrace{F(\mathbf{x}^k, \boldsymbol{\Xi}_k) + \langle \tilde{\nabla}_k, \mathbf{x} - \mathbf{x}^k \rangle}_{F_{\mathbf{x}^k}(\mathbf{x}): \text{model of } F \text{ at } \mathbf{x}^k} + \frac{1}{\alpha_k} \mathcal{D}_\phi(\mathbf{x}, \mathbf{x}^k) + R(\mathbf{x}) \right\}, \quad (3)$$

where \mathcal{D}_ϕ is the Bregman distance induced by the kernel function ϕ . To illustrate the advantage of the Bregman proximity model, we present a toy example. Consider the objective function $F(x) = x^4$, which does not admit a globally Lipschitz continuous gradient. We compare the performance of the upper quadratic approximation model (2) and Bregman proximity model (3). As shown in Figure (1)(a), the Bregman proximity model (3) ($F_2(x)$) with the kernel function $\phi(x) = \frac{1}{2}x^2 + \frac{1}{4}x^4$ can provide a closer approximation to $F(x)$ than the upper quadratic approximation model 2 ($F_1(x)$), as the yellow curve stays closer to the curve of the objective function $F(x) = x^4$. This improved approximation enables the Bregman gradient method to generate x^{k+1} that makes more substantial progress toward the optimal solution ($x^* = 0$) compared to the gradient descent method, as shown in Figure 1(b).

While several stochastic extensions of Bregman methods based on smooth adaptivity assumptions have been developed recently, the existing literature primarily focuses on stochastic convex problems (Dragomir et al., 2021b; Hanzely and Richtárik, 2021; Lu, 2019). Convergence analyses for Bregman methods in nonconvex settings (Latafat et al., 2022; Wang and Han, 2023) rely heavily on variance reduction techniques and finite sum structure.

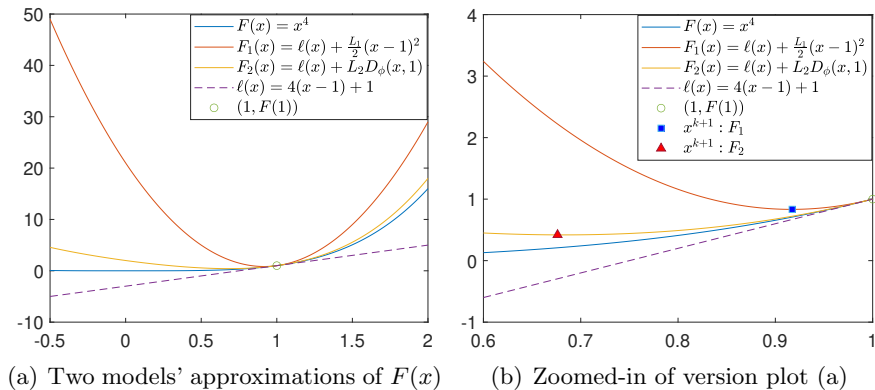


Figure 1: For function $F(x) = x^4$, which does not admit a globally Lipschitz continuous gradient. We restrict the feasible set to $[-0.5, 2]$. Consider the models (2) and (3) of F at $x^k = 1$. The Lipschitz smooth constant of F with respect to the kernel $\phi(x) = \frac{1}{2}x^2$ is 48. The smooth adaptivity constant of F with respect to the kernel $\phi(x) = \frac{1}{2}x^2 + \frac{1}{4}x^4$ is 4. The figure in (b) is a zoomed-in version of the plot in (a) for the range $[0.6, 1]$. The unique minimum of $F(x)$ is at $x = 0$.

However, these approaches are either memory-intensive or require periodic computation of the full gradient, making them impractical for large-scale problems such as deep neural networks as demonstrated in (Defazio and Bottou, 2019).

Beyond the optimization literature, recent works in learning theory (Azizan et al., 2019; Li et al., 2021; Sun et al., 2022, 2023) have focused on specific scenarios, such as highly overparameterized models or linear models for classification and regression tasks, examining the implicit bias of Bregman-type methods. For instance, (Azizan et al., 2019) demonstrates that in overparameterized models, mirror descent iterates converge to the global minimum closest to the initial point in terms of Bregman divergence, provided the initial point is near the minimum set. Other works (Li et al., 2021; Sun et al., 2022, 2023) explore convergence towards the direction that maximizes the margin in classification problems. Our work, however, focuses on the optimization properties of stochastic Bregman proximal gradient (SBPG) in general nonconvex optimization problems, which presents a different perspective from the implicit bias literature. As a result, we will not review the broader literature on implicit bias in this context.

As we can see, stochastic Bregman methods have not been fully explored in the general large-scale nonconvex problems such as training neural networks, and rigorous numerical evaluations of their performance are limited. Additionally, the current literature gives insufficient attention to the robustness of stochastic Bregman methods, particularly regarding the selection of stepsizes and initial points—factors that can have a substantial impact on their effectiveness in large-scale applications.

In this paper, we consider stochastic Bregman proximal gradient methods (SBPG) for nonconvex problems with application to the training of deep neural networks. We establish the convergence result of a vanilla SBPG without Lipschitz smoothness assumption for

nonconvex problems. Moreover, we propose a momentum-based variant, denoted as MSBPG, and prove that it offers improved convergence properties compared to vanilla SBPG, particularly through a relaxed requirement on the mini-batch size. Both methods exhibit sample complexity of $\mathcal{O}(\epsilon^{-4})$, which matches the optimal bound for stochastic first order methods (Arjevani et al., 2023). We apply MSBPG to the training of deep neural networks with a polynomial kernel function, which ensures the smooth adaptivity of the loss function. Through an implicit reformulation, we observe that MSBPG enhances the robustness of neural network training by mitigating the gradient explosion phenomenon. For numerical illustrations, we conduct experiments on quadratic inverse problems (QIP) and testify vanilla SBPG’s robustness to stepsize selection and initial point scaling. We also conduct extensive experiments on training deep neural networks for image classification and LSTMs for language modeling by employing MSBPG, which is well-suited for solving large-scale problems. Experimental results on representative benchmarks show that MSBPG converges to stationary points and, in some cases, nearly achieves the global minimum (i.e., training loss approaches zero). Moreover, MSBPG achieves comparable generalization performance to widely-used optimization algorithms, such as SGD (Robbins and Monro, 1951), Adam (Kingma and Ba, 2014), and AdamW (Loshchilov and Hutter, 2017), and in some instances, even outperforms these methods. The polynomial kernel function employed contributes to improved algorithmic stability compared to standard SGD, which may partly explain the good generalization performance of MSBPG observed in our experiments (see Appendix B for details). As the primary focus of this paper is on the optimization properties of Bregman-type methods, we only provide a preliminary analysis of algorithmic stability in the appendix. Furthermore, compared with standard SGD, SBPG/MSBPG is more robust to large stepsize and initial point scaling, which are the common reasons behind gradient explosion.

To summarize, our contributions are as follows:

1. **Development of SBPG for General Nonconvex Composite Problems:** We investigate Stochastic Bregman Proximal Gradient (SBPG) method to solve nonconvex problems without finite-sum structure, which employs Bregman distance to handle the non-Lipschitz gradient continuity. we establish its convergence properties along with optimal sample complexity of $\mathcal{O}(\epsilon^{-4})$. Furthermore, we propose a momentum-based variant, MSBPG, which improves the convergence property by relaxing the mini-batch size requirements, which is more suitable for large-scale problems. To our knowledge, this is the first integration of momentum techniques into a stochastic Bregman proximal gradient framework for nonconvex problems.
2. **Tailored MSBPG for Deep Neural Networks with Polynomial Kernel:** We apply MSBPG to training deep neural networks (DNN), which leverages on a suitable polynomial kernel function to ensure that the DNN’s loss function is smooth adaptable with respect to the designed kernel function. Through an implicit reformulation, we observe that MSBPG is more robust than the traditional SGD in terms of stepsize selection and initialization. We highlight that MSBPG is a theoretically derived method that is able to ease the difficulty of selecting stepsize, mitigate gradient explosion, and maintain good generalization performance simultaneously. This distinguishes MSBPG from many existing techniques that rely on intuition and empirical observations.

3. **Empirical Evaluation of SBPG/MSBPG:** We demonstrate the efficiency and robustness of SBPG/MSBPG across various applications, including sparse quadratic inverse problems and large-scale deep neural networks. In the quadratic inverse problem, SBPG shows greater robustness to both stepsize and initial point selection. When training deep neural networks, MSBPG achieves comparable generalization performance to commonly used optimizers such as SGD, Adam, and AdamW, and in many cases, even outperforms them. Additionally, our method demonstrates robustness to stepsize selection and initialization. These results highlight the potential of MSBPG as a powerful tool for optimizing complex and large-scale deep neural networks, thus offering a promising direction for future research in this area.

The remainder of this paper is organized as follows. In Section 2, we present notation, some related preliminaries, and our problem setting. In Section 3, we first describe SBPG and establish its convergence results. Then, we propose a momentum-based SBPG (MSBPG) and prove its improved convergence properties. In Section 4, we adapt SBPG/MSBPG to the training of deep neural networks and analyze its capacity in mitigating gradient explosion. In Section 5, we present numerical experiments that demonstrate the efficiency and robustness of vanilla SBPG on quadratic inverse problems and MSBPG on training deep neural networks. Finally, we give some concluding remarks in Section 6. Additional supplementary materials are provided in the Appendix.

2. Preliminaries and Problem setting

In this paper, vectors are represented using boldface letters like \mathbf{v} , while scalars are represented using normal font. Given a proper, lower-semicontinuous function $F : \mathbb{R}^d \rightarrow \bar{\mathbb{R}} := [-\infty, \infty]$, $\text{dom } F = \{\mathbf{x} : F(\mathbf{x}) < \infty\}$. The Fenchel conjugate function of F is defined as $F^*(\mathbf{y}) = \sup\{\langle \mathbf{x}, \mathbf{y} \rangle - F(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d\}$. Given a set $\mathcal{S} \subset \mathbb{R}^d$, $\bar{\mathcal{S}}$ denotes its closure, $\text{int } \mathcal{S}$ denotes the set of interior points. A function is of class $\mathcal{C}^k(\mathcal{S})$ if it is k times differentiable and the k -th derivative is continuous on \mathcal{S} . We say that F is level bounded if the set $\{\mathbf{x} : F(\mathbf{x}) < \alpha\}$ is bounded for any real number α . Given a matrix A , $\text{Vec}(A)$ denotes the vectorization of A by column order. $\text{Mat}(\cdot)$ is the inverse operation of $\text{Vec}(\cdot)$, which reshapes a vector back into its original matrix form. Define the operator $\text{Diag}(\cdot)$ to map a vector into a diagonal matrix with diagonal elements equal to the corresponding entries of the vector. The Hadamard product is represented by the symbol \circ . If we use the notation $\|\cdot\|$ without any additional explanation, we assume that it refers to the Euclidean vector norm for vectors and the Frobenius matrix norm for matrices.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. Given a random variable $\boldsymbol{\xi}$ and a σ -algebra \mathcal{F} , we write $\boldsymbol{\xi} \triangleleft \mathcal{F}$ if $\boldsymbol{\xi}$ is measurable over \mathcal{F} . Let $\{\boldsymbol{\xi}_k\}_{k \geq 0}$ be a stochastic process, and $\{\mathcal{F}_k\}_{k \geq 0}$ be a filtration, where \mathcal{F}_k defined by a σ -algebra $\mathcal{F}_k := \sigma(\boldsymbol{\xi}_0, \dots, \boldsymbol{\xi}_{k-1})$ on Ω . The conditional expectation is denoted by $\mathbb{E}[\cdot | \mathcal{F}_k]$. For simplicity, we use the notation $\mathbb{E}[\cdot]$ to denote $\mathbb{E}[\cdot | \mathcal{F}_\infty]$. The sequence $\{\mathbf{x}^k\}_{k \geq 0}$ generated by our proposed method is adapted to the filtration $\{\mathcal{F}_k\}_{k \geq 0}$, i.e. $\mathbf{x}^k \triangleleft \mathcal{F}_k$, for all $k \geq 0$. The notation $\tilde{\nabla}_k$ represents an estimator of the exact gradient $\nabla F(\mathbf{x}^k)$, which satisfies $\tilde{\nabla}_k \triangleleft \mathcal{F}_{k+1}$. This estimator is applicable to both the vanilla and momentum cases. The stochastic error is denoted by $\boldsymbol{\varepsilon}_k = \nabla F(\mathbf{x}^k) - \tilde{\nabla}_k$. The unbiasedness of the stochastic error $\boldsymbol{\varepsilon}_k$ is assumed throughout this paper, i.e., $\mathbb{E}[\boldsymbol{\varepsilon}_k | \mathcal{F}_k] = 0$.

The following supermartingale convergence lemma is a fundamental tool in the analysis of stochastic algorithms.

Lemma 1 (*Robbins and Monro, 1951*) *Let $\{y_k\}$, $\{u_k\}$, $\{a_k\}$ and $\{b_k\}$ be non-negative adapted processes with respect to the filtration $\{\mathcal{F}_k\}$ such that $\sum_{k=0}^{\infty} a_k < \infty$, $\sum_{k=0}^{\infty} b_k < \infty$, and for all k , $\mathbb{E}[y_{k+1} \mid \mathcal{F}_k] \leq (1 + a_k)y_k - u_k + b_k$ almost surely. Then, $\{y_k\}$ converges almost surely to a non-negative finite random variable and $\sum_{k=0}^{\infty} u_k < \infty$ almost surely.*

2.1 Smooth adaptable functions

In this subsection, we introduce the concept of smooth adaptivity, initially proposed by Bolte et al. (2018). This concept extends the idea of relative smoothness for convex functions, first introduced in Bauschke et al. (2017); Lu et al. (2018), and has since been applied in various contexts (Hanzely et al., 2021; Yang and Toh, 2021, 2022; Yang et al., 2024). We first give the definitions of kernel function and Bregman distance.

Definition 2 (*Kernel function and Bregman distance*). *Let \mathcal{S} be a nonempty, convex and open subset of \mathbb{R}^d . A function $\phi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is called a kernel function associated with \mathcal{S} if it satisfies the following two conditions:*

1. ϕ is proper, lower-semicontinuous and convex, with $\text{dom } \phi \subset \bar{\mathcal{S}}$, $\text{dom } \partial\phi = \mathcal{S}$.
2. $\phi \in \mathcal{C}^1(\mathcal{S})$ and $\text{int } \text{dom } \phi = \mathcal{S}$.

Denote the class of kernel function associated with \mathcal{S} by $\mathcal{M}(\mathcal{S})$. Given $\phi \in \mathcal{M}(\mathcal{S})$, the Bregman distance (Bregman, 1967) generated by ϕ is defined as $\mathcal{D}_\phi(\mathbf{x}, \mathbf{y}) : \text{dom } \phi \times \text{int } \text{dom } \phi \rightarrow [0, +\infty)$, where

$$\mathcal{D}_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla\phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle.$$

Bregman distance measures the difference between the value of ϕ at \mathbf{x} and its linear approximation at \mathbf{y} based on the gradient of ϕ at \mathbf{y} . Some basic properties of Bregman distance can be found in (Chen and Teboulle, 1993; Teboulle, 2018). Some kernel functions commonly used in optimization are $\frac{1}{2}\|\mathbf{x}\|^2$, $\frac{1}{2}\|\mathbf{x}\|^2 + \frac{\alpha}{4}\|\mathbf{x}\|^4$, $-\sum_{i=1}^d \log \mathbf{x}_i$ and $\sum_{i=1}^d \mathbf{x}_i \log \mathbf{x}_i$, where $\frac{1}{2}\|\mathbf{x}\|^2 \in \mathcal{M}(\mathbb{R}^d)$ recovers the classical half squared Euclidean distance. The kernel function $\frac{1}{2}\|\mathbf{x}\|^2 + \frac{\alpha}{4}\|\mathbf{x}\|^4 \in \mathcal{M}(\mathbb{R}^d)$ has found applications in various problems, such as quadratic inverse problems, non-negative matrix factorization, and low-rank minimization (Bolte et al., 2018; Dragomir et al., 2021a). The entropy function $\sum_{i=1}^d \mathbf{x}_i \log \mathbf{x}_i \in \mathcal{M}(\mathbb{R}_{++}^d)$ is commonly used in applications that involve probability constraints, where the resulting Bregman distance is known as the Kullback–Leibler (KL) divergence. Throughout the paper we will focus on the following pair of functions (f, ϕ) satisfying smooth adaptivity condition. We introduce this concept in the following definition:

Definition 3 (*Smooth adaptivity*). *Given a kernel function $\phi \in \mathcal{M}(\mathcal{S})$, a proper lower-semicontinuous function $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ with $\text{dom } f \supset \text{dom } \phi$ that is \mathcal{C}^1 on \mathcal{S} . f is L -smooth adaptable with respect to ϕ if there exists $L > 0$, such that $L\phi + f$ and $L\phi - f$ are convex on \mathcal{S} .*

Alternative definition of smooth adaptivity is the two-side descent lemma (Bolte et al., 2018, Lemma 2.1). When both f and ϕ belong to $\mathcal{C}^2(\mathcal{S})$, we can verify their smooth adaptivity by comparing the Hessians of f and ϕ .

Lemma 4 *f is L -smooth adaptable with respect to $\phi \in \mathcal{M}(\mathcal{S})$, if and only if*

$$|f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle| \leq L\mathcal{D}_\phi(\mathbf{x}, \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \text{int dom } \phi.$$

Moreover, when both f and ϕ belong to $\mathcal{C}^2(\text{int dom } \phi)$, then the above is equivalent to

$$\exists L > 0, L\nabla^2\phi(\mathbf{x}) \pm \nabla^2 f(\mathbf{x}) \succeq 0, \quad \text{for all } \mathbf{x} \in \text{int dom } \phi.$$

The following four-point identity is frequently employed in our proofs, and can be easily verified.

Lemma 5 (*Four points identity*) *Given points $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}$ and any convex function ϕ which is differentiable at \mathbf{a} and \mathbf{b} , then*

$$\langle \nabla\phi(\mathbf{a}) - \nabla\phi(\mathbf{b}), \mathbf{c} - \mathbf{d} \rangle = \mathcal{D}_\phi(\mathbf{c}, \mathbf{b}) + \mathcal{D}_\phi(\mathbf{d}, \mathbf{a}) - \mathcal{D}_\phi(\mathbf{c}, \mathbf{a}) - \mathcal{D}_\phi(\mathbf{d}, \mathbf{b}).$$

2.2 Bregman Proximal Mapping

Throughout this paper, we make the following basic assumptions.

Assumption 1 (*Basic requirements*). *In problem (1):*

- A1. *F is a proper lower-semicontinuous function with $\text{dom } \phi \subset \text{dom } F$, and it is \mathcal{C}^1 on $\text{int } C$.*
- A2. *The Legendre kernel (Definition 6) $\phi \in \mathcal{M}(C)$ is μ -strongly convex for some $\mu > 0$. $F(\cdot)$ is L_F -smooth adaptable with respect to ϕ .*
- A3. *R is a proper, lower-semicontinuous and convex function with $\text{dom } R \cap \text{int } C \neq \emptyset$.*
- A4. *$\inf_{\mathbf{x} \in \overline{C}} \{\Phi(\mathbf{x})\} > -\infty$.*

Assumption 1 is a standard requirement for Bregman-type methods and is usually satisfied in practice. It ensures the well-definedness of Bregman-type methods, as shown in (Bolte et al., 2018; Latafat et al., 2022). We also recall the definition of the Legendre function in (Latafat et al., 2022), which makes additional supercoercive conditions on the concept in (Rockafellar, 1997).

Definition 6 (*Legendre kernel*). *Let $\phi : \overline{C} \rightarrow (-\infty, \infty]$ be a proper lower-semicontinuous convex function. It is called essentially smooth if $\text{int dom } \phi$ is nonempty and ϕ is differentiable on $\text{int dom } \phi$, moreover $\lim_{k \rightarrow \infty} \|\nabla\phi(\mathbf{x}^k)\| = \infty$ whenever $\{\mathbf{x}^k\}_{k \in \mathbb{N}}$ converges to a boundary point of $\text{dom } \phi$. The function ϕ is called Legendre function if it is essentially smooth, strictly convex on $\text{int dom } \phi$ and supercoercive, i.e. $\lim_{\|\mathbf{x}\| \rightarrow \infty} \frac{\phi(\mathbf{x})}{\|\mathbf{x}\|} = \infty$.*

Definition 7 *Given a nonempty convex open set C , a proper lower-semicontinuous convex function R and a Legendre kernel function $\phi \in \mathcal{M}(C)$, $\mathbf{x} \in \text{int dom } \phi$, we denote the Bregman proximal mapping by $\text{Prox}_R^\phi := (\nabla\phi + \partial R)^{-1}\nabla\phi$, which is equivalent to*

$$\text{Prox}_R^\phi(\mathbf{x}) = \underset{\mathbf{u} \in \overline{C}}{\text{argmin}} \{R(\mathbf{u}) + \mathcal{D}_\phi(\mathbf{u}, \mathbf{x})\}. \quad (4)$$

Note that the objective function of (4) is strictly convex on $\text{dom } \phi \cap \text{dom } R$, therefore (4) has at most one solution. To ensure that (4) is well-defined, the following result claims that $\text{Prox}_{\alpha R}^{\phi}(\mathbf{x})$ is well-defined for any $\alpha > 0$, and moreover $\text{Prox}_{\alpha R}^{\phi}(\mathbf{x}) \in \text{int } \text{dom } \phi$ under standard assumptions. The proof can be found in Appendix A.

Lemma 8 *Suppose Assumption 1 holds. Then (4) has a unique solution. Moreover, the solution $\text{Prox}_{\alpha R}^{\phi}(\mathbf{x}) \in C$.*

The following proposition for Bregman proximal mapping generalizes the nonexpansive property of the classical proximal mapping (in the case $\phi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$). This property is commonly used in convergence proofs. The proof of the following proposition can be found in Appendix A.

Proposition 9 *Suppose Assumption 1 holds. Let $\mathbf{x}_i^+ := \text{Prox}_R^{\phi}(\nabla\phi^*(\mathbf{x}_i))$, $i = 1, 2$. Then $\|\mathbf{x}_1^+ - \mathbf{x}_2^+\| \leq \frac{1}{\mu}\|\mathbf{x}_1 - \mathbf{x}_2\|$.*

In this paper, we make the assumption that R and ϕ are simple enough so that (4) either has a closed-form solution or admits an efficient subroutine to solve it. Using the definition of the Bregman proximal mapping, we can then define the Bregman gradient mapping associated with (1). This mapping measures the solution accuracy of the methods we propose. Note that ϕ is a Legendre kernel, which implies that $\phi^* \in \mathcal{C}^1(\mathbb{R}^d)$ is strictly convex and $(\nabla\phi)^{-1} = \nabla\phi^*$ (Rockafellar, 1997, Corollary 13.3.1, Theorem 26.5). Therefore, the following concept is well-defined.

Definition 10 (Bregman Gradient Mapping) *Given $\alpha > 0$, a nonempty convex open set C and a Legendre kernel function $\phi \in \mathcal{M}(C)$, the Bregman gradient mapping associated with (1) is defined as follows*

$$\mathcal{G}_{\alpha}(\mathbf{x}) = \frac{\mathbf{x} - \text{Prox}_{\alpha R}^{\phi}(\nabla\phi^*(\nabla\phi(\mathbf{x}) - \alpha\nabla F(\mathbf{x})))}{\alpha}.$$

To simplify notation, we use $\mathcal{G}(\mathbf{x})$ to denote $\mathcal{G}_1(\mathbf{x})$ when $\alpha = 1$.

When the kernel function $\phi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$, the resulting Bregman Gradient Mapping becomes equivalent to the classical Gradient Mapping (Nesterov, 2003, 2005), which measures the solution's accuracy for proximal gradient methods.

Definition 11 (Limiting subdifferential (Rockafellar and Wets, 1998, Definition 8.3)) *Consider a function $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ and a point \mathbf{x} , the regular subdifferential is defined as*

$$\hat{\partial}f(\mathbf{x}) = \{\mathbf{v} : f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{v}, \mathbf{y} - \mathbf{x} \rangle + o(\|\mathbf{y} - \mathbf{x}\|)\}.$$

The limiting subdifferential is defined as

$$\partial f(\mathbf{x}) = \{\mathbf{v} : \mathbf{x}_n \rightarrow \mathbf{x}, f(\mathbf{x}_n) \rightarrow f(\mathbf{x}), \mathbf{v}_n \in \hat{\partial}f(\mathbf{x}_n), \text{ and } \mathbf{v}_n \rightarrow \mathbf{v}\}.$$

Now, we restrict our attention on the case $C = \mathbb{R}^d$. By Fermat's rule (Rockafellar and Wets, 1998, Theorem 10.1), the set of critical point of Φ is given by

$$\text{crit } \Phi = \left\{ \mathbf{x} \in \mathbb{R}^d : 0 \in \partial\Phi(\mathbf{x}) \equiv \nabla F(\mathbf{x}) + \partial R(\mathbf{x}) \right\}.$$

The Bregman Gradient Mapping can also be used to evaluate the solution accuracy for Bregman methods. Let $\mathbf{x}^+ = \text{Prox}_{\alpha R}^{\phi}(\nabla\phi^*(\nabla\phi(\mathbf{x}) - \alpha\nabla F(\mathbf{x})))$. From Definition 10 and equation (4), it can be easily verified by definition that $0 \in \partial\Phi(\mathbf{x}) \Leftrightarrow 0 = \mathcal{G}_{\alpha}(\mathbf{x})$. Hence, $0 \in \partial\Phi(\mathbf{x}^+)$ for any $\alpha > 0$. The proof of this result is omitted for brevity. Furthermore, if $\nabla\phi$ is L_{ϕ} -Lipschitz continuous, then the following proposition holds, implying that $\|\mathcal{G}_{\alpha}(\mathbf{x})\|$ can be used as a reasonable criterion to measure the accuracy of \mathbf{x} .

Proposition 12 *Suppose Assumption 1 holds and that $\nabla\phi$ is L_{ϕ} Lipschitz continuous. Then, we have the following inequality:*

$$\text{dist}(0, \partial\Phi(\mathbf{x}^+)) \leq (1 + \alpha L_F)L_{\phi}\|\mathcal{G}_{\alpha}(\mathbf{x})\|.$$

We also define the stochastic counterpart of Definition 10, which is commonly utilized to evaluate the accuracy of solutions for nonconvex stochastic proximal gradient methods, as discussed in (Ghadimi et al., 2016).

Definition 13 (*Stochastic Bregman Gradient Mapping*). *Given $\alpha > 0$ a nonempty convex open set C and a Legendre kernel function $\phi \in \mathcal{M}(C)$, the stochastic Bregman gradient mapping associated with (1) is defined as follows*

$$\tilde{\mathcal{G}}_{\alpha}(\mathbf{x}) := \frac{\mathbf{x} - \text{Prox}_{\alpha R}^{\phi}\left(\nabla\phi^*\left(\nabla\phi(\mathbf{x}) - \alpha\tilde{\nabla}\right)\right)}{\alpha}, \text{ where } \tilde{\nabla} \text{ is an estimator of } \nabla F(\mathbf{x}).$$

3. Stochastic Bregman Proximal Gradient Method

In this section, we will study the Stochastic Bregman Proximal Gradient method (SBPG) with the following update scheme:

$$\mathbf{x}^{k+1} = \underset{\mathbf{x} \in \bar{C}}{\text{argmin}} R(\mathbf{x}) + \langle \tilde{\nabla}_k, \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{\alpha_k} \mathcal{D}_{\phi}(\mathbf{x}, \mathbf{x}^k). \quad (\text{SBPG})$$

We call the above method as "vanilla" SBPG in this section, meaning that the method we study is a basic version without any additional techniques such as variance reduction, momentum, etc., except for the use of mini-batches. In this case, we suppose the following assumptions.

Assumption 2 (*Noise requirement*). *The estimator satisfies the following two conditions:*

$$\mathbb{E}[\tilde{\nabla}_k | \mathcal{F}_k] = \nabla F(\mathbf{x}_k) \quad \text{and} \quad \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] \leq \frac{\sigma^2}{m_k},$$

where m_k is the size of the mini-batch in the k -th iteration.

Note that we do not assume a finite-sum structure for $F(\mathbf{x})$ in this section. The solution of (SBPG) can be written in the form of the Bregman proximal mapping. This is stated in the following proposition.

Proposition 14 *Suppose Assumption 1 holds. Then the solution of (SBPG) can be written as the following Bregman proximal mapping:*

$$\mathbf{x}^{k+1} = \text{Prox}_{\alpha_k R}^{\phi}\left(\nabla\phi^*\left(\nabla\phi(\mathbf{x}^k) - \alpha_k\tilde{\nabla}_k\right)\right).$$

Proof From the optimality condition of the main subproblem (SBPG), we have

$$0 \in \partial R(\mathbf{x}^{k+1}) + \tilde{\nabla}_k + \frac{1}{\alpha_k} \left(\nabla \phi(\mathbf{x}^{k+1}) - \nabla \phi(\mathbf{x}^k) \right).$$

Let $\mathbf{u}^{k+1} := \text{Prox}_{\alpha_k R}^{\phi} \left(\nabla \phi^* \left(\nabla \phi(\mathbf{x}^k) - \alpha_k \tilde{\nabla}_k \right) \right)$. From the definition of Bregman proximal mapping, we have

$$\mathbf{u}^{k+1} = \arg \min_{\mathbf{u}} \left\{ \alpha_k R(\mathbf{u}) + \mathcal{D}_{\phi} \left(\mathbf{u}, \nabla \phi^* \left(\nabla \phi(\mathbf{x}^k) - \alpha_k \tilde{\nabla}_k \right) \right) \right\},$$

which is equivalent to

$$0 \in \alpha_k \partial R(\mathbf{u}^{k+1}) + \nabla \phi(\mathbf{u}^{k+1}) - \nabla \phi(\nabla \phi^*(\nabla \phi(\mathbf{x}^k) - \alpha_k \tilde{\nabla}_k)).$$

Note that the function ϕ^* is the Fenchel conjugate of the Legendre kernel ϕ , which implies that $\nabla \phi(\nabla \phi^*(\mathbf{w})) = \mathbf{w}$ for all $\mathbf{w} \in \mathbb{R}^d$, as stated in (Rockafellar, 1997, Corollary 13.3.1, Theorem 26.5). Furthermore, since the objective function in (SBPG) is strictly convex, there exists a unique solution to the inclusion above. By comparing the two inclusions, we can conclude that $\mathbf{u}^{k+1} = \mathbf{x}^{k+1}$. \blacksquare

Based on Proposition 14 and definition of the definition of $\tilde{\mathcal{G}}_{\alpha}(\mathbf{x})$, we can easily observe that $\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)$. We can derive the following proposition, which bounds the difference between $\mathcal{G}_{\alpha}(\mathbf{x})$ and $\tilde{\mathcal{G}}_{\alpha}(\mathbf{x})$ directly from Proposition 9. The proof is omitted for brevity.

Proposition 15 *Suppose Assumption 1 holds. At the k -th step, we have the estimation:*

$$\|\mathcal{G}_{\alpha_k}(\mathbf{x}^k) - \tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\| \leq \frac{1}{\mu} \|\nabla F(\mathbf{x}^k) - \tilde{\nabla}_k\| = \frac{\|\boldsymbol{\varepsilon}_k\|}{\mu},$$

where $\boldsymbol{\varepsilon}_k = \nabla F(\mathbf{x}^k) - \tilde{\nabla}_k$.

Before presenting the main convergence result, we state the following one-step descent lemma below.

Lemma 16 *Suppose Assumption 1 holds. The sequence generated by SBPG satisfies the following condition:*

$$\Phi(\mathbf{x}^{k+1}) \leq \Phi(\mathbf{x}^k) - \frac{1}{\alpha_k} \mathcal{D}_{\phi}(\mathbf{x}^k, \mathbf{x}^{k+1}) - \left(\frac{1}{\alpha_k} - L_F \right) \mathcal{D}_{\phi}(\mathbf{x}^{k+1}, \mathbf{x}^k) + \langle \boldsymbol{\varepsilon}_k, \mathbf{x}^{k+1} - \mathbf{x}^k \rangle.$$

Proof By the optimality condition of (SBPG), we obtain that

$$0 \in \partial R(\mathbf{x}^{k+1}) + \tilde{\nabla}_k + \frac{1}{\alpha_k} \left(\nabla \phi(\mathbf{x}^{k+1}) - \nabla \phi(\mathbf{x}^k) \right).$$

Appealing to the convexity of R , we have

$$R(\mathbf{x}) - R(\mathbf{x}^{k+1}) \geq \langle -\tilde{\nabla}_k - \frac{1}{\alpha_k} \left(\nabla \phi(\mathbf{x}^{k+1}) - \nabla \phi(\mathbf{x}^k) \right), \mathbf{x} - \mathbf{x}^{k+1} \rangle.$$

By the four points identity and the definition of ε_k , we get

$$R(\mathbf{x}) - R(\mathbf{x}^{k+1}) \geq \frac{1}{\alpha_k} \left[\mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k) + \mathcal{D}_\phi(\mathbf{x}, \mathbf{x}^{k+1}) - \mathcal{D}_\phi(\mathbf{x}, \mathbf{x}^k) \right] - \langle \nabla F(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^{k+1} \rangle + \langle \varepsilon_k, \mathbf{x} - \mathbf{x}^{k+1} \rangle.$$

Set $\mathbf{x} = \mathbf{x}^k$ in the above inequality, we have the following inequality:

$$R(\mathbf{x}^k) - R(\mathbf{x}^{k+1}) \geq \frac{1}{\alpha_k} \left[\mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k) + \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) \right] - \langle \nabla F(\mathbf{x}^k), \mathbf{x}^k - \mathbf{x}^{k+1} \rangle + \langle \varepsilon_k, \mathbf{x}^k - \mathbf{x}^{k+1} \rangle.$$

By the smooth adaptivity of F , we have

$$F(\mathbf{x}^{k+1}) \leq F(\mathbf{x}^k) + \langle \nabla F(\mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + L_F \mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k).$$

Combining the above two inequalities above, we complete the proof. \blacksquare

3.1 Convergence analysis of SBPG

In this subsection, we establish the convergence results for SBPG, which is an extension of the convergence result in (Ghadimi et al., 2016), in which the lassical Lipschitz gradient assumption is required. In many literature, the bounded sequence assumption is often required in the convergence analysis of stochastic algorithms. However, in this section, we relax this assumption and prove that under a certain condition, the sequence generated by (SBPG) is bounded almost surely. We need the following result to bound the stochastic error term $\langle \varepsilon_k, \mathbf{x}^{k+1} - \mathbf{x}^k \rangle$ in Lemma 16.

Lemma 17 *Suppose Assumption 1, 2 hold. We have the following estimation of the error term:*

$$\mathbb{E} \left[\langle \varepsilon_k, \mathbf{x}^{k+1} - \mathbf{x}^k \rangle \right] \leq \frac{\alpha_k}{\mu} \mathbb{E}[\|\varepsilon_k\|^2] \leq \frac{\alpha_k \sigma^2}{\mu m_k}.$$

Proof Define $\bar{\mathbf{x}}^{k+1} := \text{Prox}_{\alpha_k R}^\phi(\nabla \phi^*(\nabla \phi(\mathbf{x}^k) - \alpha_k \nabla F(\mathbf{x}^k)))$. By Proposition 14 and the optimality condition for $\bar{\mathbf{x}}^{k+1}$, we have

$$0 \in \partial R(\bar{\mathbf{x}}^{k+1}) + \nabla F(\mathbf{x}^k) + \frac{1}{\alpha_k} (\nabla \phi(\bar{\mathbf{x}}^{k+1}) - \nabla \phi(\mathbf{x}^k)).$$

Similarly,

$$0 \in \partial R(\mathbf{x}^{k+1}) + \tilde{\nabla}_k + \frac{1}{\alpha_k} (\nabla \phi(\mathbf{x}^{k+1}) - \nabla \phi(\mathbf{x}^k)).$$

By the monotonicity of ∂R and Lemma 5, we have

$$\langle \bar{\mathbf{x}}^{k+1} - \mathbf{x}^{k+1}, -\varepsilon_k - \frac{1}{\alpha_k} (\nabla \phi(\bar{\mathbf{x}}^{k+1}) - \nabla \phi(\mathbf{x}^{k+1})) \rangle \geq 0.$$

Therefore,

$$\langle \mathbf{x}^{k+1} - \bar{\mathbf{x}}^{k+1}, \varepsilon_k \rangle \geq \langle \bar{\mathbf{x}}^{k+1} - \mathbf{x}^{k+1}, \frac{1}{\alpha_k} (\nabla \phi(\bar{\mathbf{x}}^{k+1}) - \nabla \phi(\mathbf{x}^{k+1})) \rangle \geq \frac{\mu}{\alpha_k} \|\bar{\mathbf{x}}^{k+1} - \mathbf{x}^{k+1}\|^2.$$

By Cauchy-Schwarz inequality, we get $\|\bar{\mathbf{x}}^{k+1} - \mathbf{x}^{k+1}\| \leq \frac{\alpha_k}{\mu} \|\boldsymbol{\varepsilon}_k\|$.

Now, we are ready to prove Lemma 17. From the definition, we know that $\bar{\mathbf{x}}^{k+1} \triangleleft \mathcal{F}_k$. Therefore, $\mathbb{E}[\langle \boldsymbol{\varepsilon}_k, \mathbf{x}^k - \bar{\mathbf{x}}^{k+1} \rangle] = \mathbb{E}[\mathbb{E}[\langle \boldsymbol{\varepsilon}_k, \mathbf{x}^k - \bar{\mathbf{x}}^{k+1} \rangle | \mathcal{F}_k]] = \mathbb{E}[\langle \mathbb{E}[\boldsymbol{\varepsilon}_k | \mathcal{F}_k], \mathbf{x}^k - \bar{\mathbf{x}}^{k+1} \rangle] = 0$, where the first equality is from the tower rule of conditional expectation, the second comes from the fact that $\mathbf{x}^k - \bar{\mathbf{x}}^{k+1} \triangleleft \mathcal{F}_k$. Hence,

$$\mathbb{E}[\langle \boldsymbol{\varepsilon}_k, \mathbf{x}^{k+1} - \mathbf{x}^k \rangle] = \mathbb{E}[\langle \boldsymbol{\varepsilon}_k, \mathbf{x}^{k+1} - \bar{\mathbf{x}}^{k+1} \rangle] - \mathbb{E}[\langle \boldsymbol{\varepsilon}_k, \mathbf{x}^k - \bar{\mathbf{x}}^{k+1} \rangle] \leq \frac{\alpha_k}{\mu} \mathbb{E}[\|\boldsymbol{\varepsilon}_k\|^2] \leq \frac{\alpha_k \sigma^2}{\mu m_k},$$

which completes the proof. \blacksquare

Lemma 18 (Bounded sequence) *Suppose Assumption 1, 2 hold. If $\sum_k \frac{\alpha_k}{m_k} < \infty$, $\sup_k \alpha_k \leq \bar{\alpha} < \frac{1}{L_F}$, then,*

1. $\sum_{k=0}^{\infty} \mathbb{E}[\mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k)] < \infty$.
2. If Φ is level bounded, then $\{\mathbf{x}^k\}_{k \geq 0}$ is bounded almost surely.

Proof By Cauchy-Young inequality, we have

$$|\langle \boldsymbol{\varepsilon}_k, \mathbf{x}^k - \mathbf{x}^{k+1} \rangle| \leq \frac{\mu}{2\alpha_k} \|\mathbf{x}^k - \mathbf{x}^{k+1}\|^2 + \frac{\alpha_k}{2\mu} \|\boldsymbol{\varepsilon}_k\|^2 \leq \frac{1}{\alpha_k} \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) + \frac{\alpha_k}{2\mu} \|\boldsymbol{\varepsilon}_k\|^2.$$

By Lemma 16, we have

$$\left(\frac{1}{\alpha_k} - L_F \right) \mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k) \leq \Phi(\mathbf{x}^k) - \Phi(\mathbf{x}^{k+1}) + \frac{\alpha_k}{2\mu} \|\boldsymbol{\varepsilon}_k\|^2. \quad (5)$$

Taking conditional expectation for both sides of (5), we get

$$\mathbb{E} \left[\left(\frac{1}{\alpha_k} - L_F \right) \mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k) | \mathcal{F}_k \right] \leq \Phi(\mathbf{x}^k) - \mathbb{E}[\Phi(\mathbf{x}^{k+1}) | \mathcal{F}_k] + \frac{\alpha_k}{2\mu} \mathbb{E}[\|\boldsymbol{\varepsilon}_k\|^2 | \mathcal{F}_k].$$

Since $\sum_{k \geq 0} \frac{\alpha_k}{2\mu} \mathbb{E}[\|\boldsymbol{\varepsilon}_k\|^2 | \mathcal{F}_k] \leq \sum_{k \geq 0} \frac{\alpha_k \sigma^2}{2\mu m_k} < \infty$, applying Theorem 1, we have that $\Phi(\mathbf{x}^k)$ converges and $\sum_{k \geq 0} \mathbb{E} \left[\left(\frac{1}{\alpha_k} - L_F \right) \mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k) | \mathcal{F}_k \right] < \infty$ almost surely. By the tower rule of conditional expectation, we have $\sum_{k=0}^{\infty} \mathbb{E}[\mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k)] < \infty$. Since $\Phi(\mathbf{x}^k)$ converges almost surely, thus $\{\Phi(\mathbf{x}^k)\}_{k \geq 0}$ is bounded almost surely. By the level boundness of Φ , we deduce that $\{\mathbf{x}^k\}_{k \geq 0}$ is bounded almost surely. \blacksquare

Now, we present our main convergence result for the vanilla SBPG, which is in the sense of expectation.

Theorem 19 (Convergence result in expectation) *Suppose Assumption 1, 2 hold, $\alpha_k < \frac{1}{L_F} \min\{1, \frac{1}{\mu}\}$. Define a random variable r with the distribution $\mathbb{P}\{r = k\} = \frac{\alpha_k}{\sum_{k=0}^{N-1} \alpha_k}$ for $k = 0, \dots, N-1$. Then,*

$$\mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_r}(\mathbf{x}^r)\|^2] \leq \frac{2\Delta_0 + 2 \sum_{k=0}^{N-1} \frac{\alpha_k \sigma^2}{\mu m_k}}{\mu \sum_{k=0}^{N-1} \alpha_k}, \quad (6)$$

where $\Delta_0 := \Phi(\mathbf{x}^0) - \Phi^*$. If $\sum_k \frac{\alpha_k}{m_k} < +\infty$ and $\sum_k \alpha_k = +\infty$, then the right hand side of (6) converges to zero. Moreover, if Φ is level bounded, then the sequence $\{\mathbf{x}^k\}_{k \geq 0}$ is bounded almost surely.

Proof Note that $\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)$ and by the strongly convexity of ϕ , Lemma 16 yields

$$\begin{aligned} \mu(\alpha_k - \frac{L_F \alpha_k^2}{2}) \|\tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\|^2 &\leq \frac{1}{\alpha_k} \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) + \left(\frac{1}{\alpha_k} - L_F\right) \mathcal{D}_\phi(\mathbf{x}^{k+1}, \mathbf{x}^k) \\ &\leq \Phi(\mathbf{x}^k) - \Phi(\mathbf{x}^{k+1}) + \langle \boldsymbol{\varepsilon}_k, \mathbf{x}^{k+1} - \mathbf{x}^k \rangle. \end{aligned}$$

Taking expectations, telescoping from $k = 0 \dots N - 1$, and using Lemma 17, we obtain

$$\sum_{k=0}^{N-1} \mu(\alpha_k - \frac{\mu L_F \alpha_k^2}{2}) \mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\|^2] \leq \Phi(\mathbf{x}^0) - \Phi(\mathbf{x}^N) + \sum_{k=0}^{N-1} \frac{\alpha_k \sigma^2}{\mu m_k}. \quad (7)$$

By utilizing the inequality $\alpha_k - \frac{\mu L_F \alpha_k^2}{2} \geq \frac{\alpha_k}{2}$, the condition $\Phi(\mathbf{x}^N) \geq \Phi^*$, and considering the definition of the random variable r , we can derive (6) from (7). \blacksquare

Remark 20 We give some remarks for Theorem 19.

1. The mini-batch size plays a crucial role in ensuring convergence, as it allows us to control the stochastic error term in Lemma 16 and provide a bound for $\mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\|^2]$ that converges to 0 as k tends to infinity. If $m_k = 1$ for all k , then the upper bound for $\mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\|^2]$ will not converge to 0, no matter how $\{\alpha_k\}$ is selected.
2. In (Ghadimi et al., 2016), a similar convergence result is established for mini-batch stochastic proximal gradient methods, but our analysis differs in a crucial aspect in that we do not assume the Lipschitz continuity of $F(\mathbf{x})$. Instead, we rely on the smooth adaptivity of $F(\mathbf{x})$, which is a more relaxed assumption. Additionally, we provide specific conditions on the stepsizes $\{\alpha_k\}$ and mini-batch sizes $\{m_k\}$ that guarantee the convergence of $\mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\|^2]$ to 0, as well as the almost sure boundedness of the sequence $\{\mathbf{x}^k\}$.
3. We now provide a specific choice of $\{\alpha_k\}$ and $\{m_k\}$ to establish a convergence rate in terms of expected stationarity. For given positive constants $c_1, c_2, \gamma \in (0, \infty)$ and $\delta \in [0, 1)$, we choose

$$\alpha_k = \frac{c_1}{(k+1)^\delta}, \quad m_k = \lceil c_2(k+1)^\gamma \rceil.$$

Under these choices, we have $\sum_{i=0}^k \alpha_i = \mathcal{O}(k^{1-\delta})$ and $\sum_{i=0}^k \frac{\alpha_i}{m_i} = \mathcal{O}(k^{1-(\delta+\gamma)})$. As a result, the RHS of (6) has the bound of $\mathcal{O}(k^{-(1-\delta)} + k^{-\gamma})$. When $\delta = 0$ and $\gamma \geq 1$, the order is $\mathcal{O}(k^{-1})$ which recovers the optimal rate of deterministic first order method. In this case, the mini-batch size m_k increases rapidly and the variance in the stochastic gradient reduces to zero quickly, so the algorithm behaves like a deterministic algorithm. However, such a favorable convergence rate comes at the expense of a heavy computational burden and high memory cost in each iterative step.

Next, we present the sample complexity of the SBPG method. We define each computation of $\nabla f(\mathbf{x}, \boldsymbol{\xi})$ for a fix $\boldsymbol{\xi}$ as an oracle evaluation.

Corollary 21 *Given an accuracy level $\epsilon > 0$ and a positive constant $\gamma > 0$, choose $\alpha_k = \alpha < \frac{1}{L_F} \min\{1, \frac{1}{\mu}\}$. Then to achieve an ϵ -stationary point in expectation, at most $\bar{N} := \left\lceil \max \left\{ \gamma^2, \frac{4\sigma^4}{\mu^4\epsilon^4} \left(\frac{2\Delta_0\mu\gamma}{\alpha\sigma^2} + \frac{1}{\gamma} \right)^2 \right\} \right\rceil$ oracle evaluations are required. In this case, we need to set $m_k = m := \lceil \min\{\gamma\sqrt{\bar{N}}, \bar{N}\} \rceil$ for all k .*

Proof Let the total number of oracle evaluations be \bar{N} . Then $N = \lfloor \frac{\bar{N}}{m} \rfloor$ and $N \geq \frac{\bar{N}}{2m}$. We have that

$$\begin{aligned} \text{RHS of (6)} &\leq \frac{4\Delta_0 m}{\alpha\mu\bar{N}} + \frac{2\sigma^2}{\mu^2 m} \leq \frac{4\Delta_0}{\alpha\mu\bar{N}} \left(\gamma\sqrt{\bar{N}} \right) + \frac{2\sigma^2}{\mu^2} \max \left\{ \frac{1}{\gamma\sqrt{\bar{N}}}, \frac{1}{\bar{N}} \right\} \\ &\leq \frac{2\sigma^2}{\mu^2\sqrt{\bar{N}}} \left(\frac{2\Delta_0\mu\gamma}{\alpha\sigma^2} + \frac{1}{\gamma} \right), \end{aligned}$$

where $\Delta_0 := \Phi(\mathbf{x}^0) - \Phi(\mathbf{x}^*)$. By the choice of \bar{N} , we have RHS of (6) $\leq \epsilon^2$. This completes the proof. \blacksquare

This oracle complexity matches the lower bound provided by (Arjevani et al., 2023), thus it cannot be improved in general.

3.2 Momentum based Stochastic Bregman Gradient Descent Method

Remark 20 and Corollary 21 suggest that a large mini-batch size m_k is necessary to achieve a small stationarity error. However, employing a large mini-batch size in each iteration can be computationally expensive and memory inefficient, especially in modern large-scale problems such as deep neural network training. In this part, we resort to the momentum technique to address this issue. Specifically, we consider using a stochastic moving average estimator (SMAE) for the true gradient given by:

$$\mathbf{v}^k = (1 - \beta_k)\mathbf{v}^{k-1} + \beta_k \tilde{\nabla}_k, \quad \text{where } \mathbb{E}[\tilde{\nabla}_k | \mathcal{F}_k] = \nabla F(\mathbf{x}^k), \quad (8)$$

where \mathbf{v}^{k-1} can be viewed as the momentum which contains the information of all historical stochastic gradients, and $\mathbb{E}[\|\tilde{\nabla}_k\|^2 | \mathcal{F}_k] \leq \frac{\sigma^2}{m_k}$. We expect that incorporating the SMAE technique can achieve a certain level of variance reduction without requiring an increase in the mini-batch size or the computation of the full gradient. In our approach, we utilize the gradient estimator \mathbf{v}^k within SBPG, and we refer to the resulting method as MSBPG. Specifically, we consider the following update scheme:

$$\mathbf{x}^{k+1} = \underset{\mathbf{x} \in \bar{C}}{\operatorname{argmin}} R(\mathbf{x}) + \langle \mathbf{v}^k, \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{\alpha_k} \mathcal{D}_\phi(\mathbf{x}, \mathbf{x}^k). \quad (9)$$

We need the following assumption that the difference of gradients of F can be bounded by the Bregman distance.

Assumption 3 *There exists $\kappa > 0$, such that $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|^2 \leq \kappa \mathcal{D}_\phi(\mathbf{x}, \mathbf{y})$ for all $\mathbf{x} \in \operatorname{dom} \phi$, $\mathbf{y} \in \operatorname{int} \operatorname{dom} \phi$.*

Remark 22 *This assumption generalizes the case of Lipschitz kernel function. If F is L_F -smooth adaptable to ϕ , it can be easily shown that if ϕ has L_ϕ -Lipschitz gradient, this assumption holds for $\kappa \geq \frac{2L_F^2 L_\phi^2}{\mu}$. In this paper, we are particularly interested in polynomial kernel functions. For functions with polynomially bounded growth rates, this assumption is not restrictive. For example, consider the one-dimensional objective function $F(x) = \frac{1}{4}x^4$ and the kernel function $\phi(x) = \frac{1}{2}x^2 + \frac{1}{8}x^8$. Then, by (Lu et al., 2018, Proposition 2.1), we know that F is smooth adaptable with respect to ϕ . Simple algebra shows that $\mathcal{D}_\phi(x, y) = \frac{1}{8}(x - y)^2(x^6 + 2x^5y + 3x^4y^2 + 4x^3y^3 + 5x^2y^4 + 6xy^5 + 7y^6 + 4)$ and $(F'(x) - F'(y))^2 = (x - y)^2(x^2 + xy + y^2)^2$. Numerical computation shows that $(x^6 + 2x^5y + 3x^4y^2 + 4x^3y^3 + 5x^2y^4 + 6xy^5 + 7y^6 + 4) - (x^2 + xy + y^2)^2 \geq 3.71$. Therefore, $(F'(x) - F'(y))^2 \leq 8\mathcal{D}_\phi(x, y)$, which holds globally for any x and y in \mathbb{R}^d .*

Next, we present a recursion lemma that allows us to estimate the accuracy of the SMAE. While similar lemmas have been proposed in the literature, such as in (Wang et al., 2017), their bounds are not directly applicable in the Bregman setting. As a result, we have developed a version of the recursion lemma that is tailored to our specific context.

Lemma 23 *The following recursion holds*

$$\mathbb{E}[\|\mathbf{v}^k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] \leq (1 - \beta_k) \|\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})\|^2 + \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] + \frac{\|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2}{\beta_k}.$$

Proof Note that $\mathbf{v}^k - \nabla F(\mathbf{x}^k) = (1 - \beta_k)(\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})) + (1 - \beta_k)(\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)) + \beta_k(\tilde{\nabla}_k - \nabla F(\mathbf{x}^k))$, and $\mathbb{E}[\tilde{\nabla}_k - \nabla F(\mathbf{x}^k) | \mathcal{F}_k] = 0$. Then we have

$$\begin{aligned} & \mathbb{E}[\|\mathbf{v}^k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] \\ &= (1 - \beta_k)^2 \|\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})\|^2 + (1 - \beta_k)^2 \|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2 + \\ & \quad \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] + 2(1 - \beta_k)^2 \langle \mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1}), \nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k) \rangle \\ &\leq (1 - \beta_k)^2 \|\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})\|^2 + (1 - \beta_k)^2 \|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2 + \\ & \quad \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] + \beta_k(1 - \beta_k) \|\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})\|^2 + \frac{(1 - \beta_k)^3}{\beta_k} \|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2 \\ &= (1 - \beta_k) \|\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})\|^2 + \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] + \frac{(1 - \beta_k)^2 \|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2}{\beta_k} \\ &\leq (1 - \beta_k) \|\mathbf{v}^{k-1} - \nabla F(\mathbf{x}^{k-1})\|^2 + \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2 | \mathcal{F}_k] + \frac{\|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2}{\beta_k}. \end{aligned}$$

This completes the proof. ■

Now we are ready to provide the convergence result for our momentum based SBPG.

Theorem 24 *Suppose Assumption 1, 2 and 3 hold. Let $\alpha_k = c\mu\beta_{k+1}$ for any $c \in (0, \frac{1}{2\sqrt{\mu\kappa}}]$. Then, it holds that*

$$\mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_r}(\mathbf{x}^r)\|^2] \leq \frac{\Delta_0 + c\|\mathbf{v}_0 - \nabla F(\mathbf{x}^0)\|^2 + \sum_{k=0}^{N-1} \frac{\alpha_k^2 \sigma^2}{c\mu^2 m_k}}{\sum_{k=0}^{N-1} \frac{\mu\alpha_k}{8}}, \quad (10)$$

where r is a random variable with distribution $\mathbb{P}\{r = k\} = \frac{\alpha_k}{\sum_{k=0}^{N-1} \alpha_k}$, for $k = 0, \dots, N - 1$.

Proof From Lemma 16 and Cauchy-Young's inequality, we have

$$\begin{aligned}\Phi(\mathbf{x}^{k+1}) &\leq \Phi(\mathbf{x}^k) - \frac{1}{\alpha_k} \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) + \frac{\mu}{4\alpha_k} \|\mathbf{x}^k - \mathbf{x}^{k+1}\|^2 + \frac{\alpha_k}{\mu} \|\boldsymbol{\varepsilon}_k\|^2 \\ &\leq \Phi(\mathbf{x}^k) - \frac{1}{2\alpha_k} \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) + \frac{\alpha_k}{\mu} \|\boldsymbol{\varepsilon}_k\|^2.\end{aligned}$$

where we have defined $\boldsymbol{\varepsilon}_k := \nabla F(\mathbf{x}^k) - \mathbf{v}^k$. Summing the above inequality over $k = 0, \dots, N-1$ and rearranging the terms, we get

$$\sum_{k=0}^{N-1} \frac{1}{2\alpha_k} \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) \leq \Phi(\mathbf{x}^0) - \Phi^* + \sum_{k=0}^{N-1} \frac{\alpha_k}{\mu} \|\boldsymbol{\varepsilon}_k\|^2.$$

By applying Lemma 23, we can obtain the following inequality:

$$\beta_k \mathbb{E}[\|\boldsymbol{\varepsilon}_{k-1}\|^2] \leq \mathbb{E}[\|\boldsymbol{\varepsilon}_{k-1}\|^2] - \mathbb{E}[\|\boldsymbol{\varepsilon}_k\|^2] + \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2] + \mathbb{E}\left[\frac{\|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2}{\beta_k}\right].$$

Hence

$$\sum_{k=0}^{N-1} \beta_{k+1} \mathbb{E}[\|\boldsymbol{\varepsilon}_k\|^2] = \sum_{k=1}^N \beta_k \mathbb{E}[\|\boldsymbol{\varepsilon}_{k-1}\|^2] \leq \|\boldsymbol{\varepsilon}_0\|^2 + \sum_{k=1}^N \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2] + \sum_{k=1}^N \mathbb{E}\left[\frac{\|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2}{\beta_k}\right].$$

Since $\frac{\alpha_k}{\mu} = c\beta_{k+1}$ for some constant c , we get the following inequality:

$$\sum_{k=0}^{N-1} \frac{1}{2\alpha_k} \mathbb{E}[\mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1})] \leq \Phi(\mathbf{x}^0) - \Phi^* + c \left(\|\boldsymbol{\varepsilon}_0\|^2 + \sum_{k=1}^N \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2] + \sum_{k=1}^N \mathbb{E}\left[\frac{\|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\|^2}{\beta_k}\right] \right).$$

By using Assumption 3, we obtain that

$$\frac{\|\nabla F(\mathbf{x}^k) - \nabla F(\mathbf{x}^{k+1})\|^2}{\beta_{k+1}} \leq \frac{\kappa}{\beta_{k+1}} \mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}).$$

Combining above two inequalities, we get

$$\sum_{k=0}^{N-1} \frac{1}{2\alpha_k} \mathbb{E}[\mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1})] \leq \Phi(\mathbf{x}^0) - \Phi^* + c \left(\|\boldsymbol{\varepsilon}_0\|^2 + \sum_{k=1}^N \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2] + \sum_{k=0}^{N-1} \frac{\kappa}{\beta_{k+1}} \mathbb{E}[\mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1})] \right).$$

Since $c \leq \frac{1}{2\sqrt{\mu\kappa}}$ and $\frac{\alpha_k}{\mu} = c\beta_{k+1}$, we can deduce that $\frac{c\kappa}{\beta_{k+1}} \leq \frac{1}{4\alpha_k}$. Using this condition, we obtain the inequality:

$$\sum_{k=0}^{N-1} \frac{1}{4\alpha_k} \mathbb{E}[\mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1})] \leq \Phi(\mathbf{x}^0) - \Phi^* + c \left(\|\boldsymbol{\varepsilon}_0\|^2 + \sum_{k=1}^N \beta_k^2 \mathbb{E}[\|\tilde{\nabla}_k - \nabla F(\mathbf{x}^k)\|^2] \right).$$

Note that $\mathcal{D}_\phi(\mathbf{x}^k, \mathbf{x}^{k+1}) \geq \frac{\mu}{2} \|\mathbf{x}^k - \mathbf{x}^{k+1}\|^2 = \frac{\mu\alpha_k^2}{2} \|\tilde{\mathcal{G}}_{\alpha_k}(\mathbf{x}^k)\|^2$ and by the definition of the random variable a , we get

$$\mathbb{E}[\|\tilde{\mathcal{G}}_{\alpha_a}(\mathbf{x}^a)\|^2] \leq \frac{\Phi^0 - \Phi^* + c\|\boldsymbol{\varepsilon}_0\|^2 + c\sum_{k=1}^N \frac{\beta_k^2 \sigma^2}{m_k}}{\sum_{k=0}^{N-1} \frac{\mu\alpha_k}{8}},$$

which completes the proof. \blacksquare

Remark 25 Now we give some remarks for Theorem 24.

1. When the sequence $\{\mathbf{x}_k\}$ is bounded, an alternative to Assumption 3 is to assume that $C = \mathbb{R}^d$ and that ϕ has a locally Lipschitz gradient, as made in (Bolte et al., 2018, Theorem 4.1) and (Latafat et al., 2022, Theorem 4.7). Under these conditions, we can conclude that there exists a compact set \mathcal{U} containing $\{\mathbf{x}_k\}$. Therefore, there exists a constant $L_{\phi, \mathcal{U}} > 0$ such that $\nabla\phi$ is Lipschitz continuous over \mathcal{U} , and we can derive that $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|^2 \leq L_F^2 L_{\phi, \mathcal{U}}^2 \|\mathbf{x} - \mathbf{y}\|^2 \leq \frac{2L_F^2 L_{\phi, \mathcal{U}}^2}{\mu} \mathcal{D}_\phi(\mathbf{x}, \mathbf{y})$ holds.
2. The stationarity error for SBPG in Theorem 19 is given by $\mathcal{O}\left(\frac{1}{\sum_{i=0}^k \alpha_i} + \frac{\sum_{i=0}^k \frac{\alpha_i}{m_i}}{\sum_{i=0}^k \alpha_i}\right)$, while for MSBPG in Theorem 24, the error is $\mathcal{O}\left(\frac{1}{\sum_{i=0}^k \alpha_i} + \frac{\sum_{i=0}^k \frac{\alpha_i^2}{m_i}}{\sum_{i=0}^k \alpha_i}\right)$. Notably, compared to SBPG, even with a small constant mini-batch size m_k , MSBPG can still achieve convergence to a zero error bound by carefully selecting the stepsize sequence $\{\alpha_k\}$. A typical stepsize condition is $\sum_{k=0}^{\infty} \alpha_k = \infty$, $\sum_{k=0}^{\infty} \alpha_k^2 < \infty$, which coincides with the classical stepsize condition ensuring a moderate decrease of the stepsize, as discussed in (Bertsekas and Tsitsiklis, 2000). Thus, by incorporating the momentum technique, we can achieve improved convergence properties, particularly in terms of relaxed mini-batch size requirements, with minimal additional computational cost. This favorable convergence property theoretically supports the application of MSBPG for large-scale problems, such as deep neural networks training, without the use of very large mini-batch size. To illustrate, we provide a specific choice of $\{\alpha_k\}$ and $\{m_k\}$ that yields a convergence rate in terms of expected stationarity: Set $m_k = 1$, $\alpha_k = \frac{c}{\sqrt{k+1}}$, the convergence rate is $\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{k}}\right)$ with logarithmic terms hidden.

Next, we present the sample complexity of MSBPG.

Corollary 26 Given an accuracy level $\epsilon > 0$, a constant $\alpha_0 > 0$, and any integer $m > 0$, set $m_k = m$. Then to achieve an ϵ -stationary point in expectation, at most $\bar{N} := \left\lceil \max \left\{ (1 + \mu^2) \alpha_0^2 L_F^2, \frac{1}{\epsilon^4} \left(\frac{16(\Delta_0 + c \|\varepsilon_0\|^2)}{\mu \alpha_0} + \frac{8\alpha_0 \sigma^2}{c \mu m} \right) \right\} \right\rceil$ oracle evaluations are required. To achieve this complexity, we can choose $\alpha_k = \alpha = \frac{\alpha_0}{\sqrt{\bar{N}}}$.

Proof Let the total number of oracle evaluations be \bar{N} . Then $N = \left\lfloor \frac{\bar{N}}{m} \right\rfloor$ and $N \geq \frac{\bar{N}}{2m}$. By choosing $\alpha_k = \alpha = \frac{\alpha_0}{\sqrt{\bar{N}}}$ and using the definition of \bar{N} , we have that $\alpha_k < \frac{1}{L_F} \min \left\{ 1, \frac{1}{\mu} \right\}$, satisfying the conditions for Theorem 24. Then we have

$$\text{RHS of (10)} \leq \frac{8(\Delta_0 + c \|\varepsilon_0\|^2)m}{\mu \alpha \bar{N}} + \frac{8\sigma^2 \alpha}{c \mu m}.$$

By the choices of \bar{N} and α , we have RHS of (24) $\leq \epsilon^2$. This completes the proof. \blacksquare

We observe that MSBPG has the same complexity order as SBPG, with both achieving $\mathcal{O}(\epsilon^{-4})$, which, as shown in (Arjevani et al., 2023), is the optimal bound and cannot be improved. However, to reach this complexity bound, SBPG requires a large mini-batch

size, whereas MSBPG allows for an arbitrary mini-batch size. This relaxed requirement for the mini-batch size makes MSBPG more practical for training deep neural networks, as it is more memory-efficient. Therefore, the convergence properties of MSBPG are improved in terms of reduced mini-batch size requirements.

As a final remark for this section, in this work, we have only established convergence to a stationary point. However, as demonstrated by (Lee et al., 2019) and (Panageas et al., 2019), the Bregman gradient method almost always converges to a local minimum when the loss function has the strict saddle property (i.e. all the saddle points of f are strict saddle points). In general, finding the global minimum of a nonconvex function is NP-hard. However, in many applications, the objective function is "benign" nonconvex, meaning that all local minima are global minima. Examples include matrix completion (Ge et al., 2016), matrix factorization (Chi et al., 2019), and phase retrieval (Sun et al., 2018), where convergence to a local minimum implies convergence to the global minimum. Additionally, when the objective function satisfies the Polyak–Łojasiewicz (PL) condition (Polyak, 1963), convergence to a stationary point also guarantees convergence to a global minimum. In the experimental section, we observe that in some instances of training neural networks, the training loss can even approach zero, indicating near convergence to a global minimum.

4. Application in training deep neural networks

In this section, we present a detailed description of MSBPG applied to training deep neural networks. Throughout this section, we assume that the optimization domain \overline{C} is the entire space \mathbb{R}^d , so that $\phi \in \mathcal{M}(\mathbb{R}^d)$ and $F \in \mathcal{C}^1(\mathbb{R}^d)$. For simplicity, we omit the explicit mention of the feasible set \mathbb{R}^d in this section. In this context, we utilize a polynomial kernel function.

The optimization problem we consider here is given by:

$$\min_{\mathbf{W}} \underbrace{\frac{1}{N} \sum_{i=1}^N \mathcal{L}(\mathcal{DNN}(\mathbf{W}, \mathbf{x}_i), y_i)}_{F(\mathbf{W})} + \lambda_2 \|\mathbf{W}\|_2^2 + \lambda_1 \|\mathbf{W}\|_1, \quad (11)$$

where $\mathcal{DNN}(\mathbf{W}, \mathbf{x})$ is the neural network function with training parameters \mathbf{W} and input data \mathbf{x} , \mathcal{L} is the loss function that measures the difference between the output of the neural network $\mathcal{DNN}(\mathbf{W}, \mathbf{x}_i)$ and the label y_i , $\lambda_2 \|\mathbf{W}\|_2^2$ is the weight decay term, and $\lambda_1 \|\mathbf{W}\|_1$ is the L_1 regularization term that is the sparsity induced operator and often used to avoid overfitting in training deep neural networks (Ng, 2004). To illustrate the neural network function $\mathcal{DNN}(\mathbf{W}, \mathbf{x})$, in the L -layer fully connected neural network, we have $\mathbf{W} = [\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L]$ and

$$\mathcal{DNN}(\mathbf{W}, \mathbf{x}) = \sigma_L(\mathbf{W}_L(\sigma_{L-1}(\mathbf{W}_{L-1}(\dots(\sigma_1(\mathbf{W}_1\mathbf{x}))\dots))), \quad (12)$$

where σ_i is the nonlinear activation function. In this paper, we focus on smooth activation functions.

At the k -th iteration, MSBPG has the following update scheme:

$$\mathbf{v}^k = (1 - \beta_k)\mathbf{v}^{k-1} + \beta_k \tilde{\nabla}_k \quad (13)$$

$$\mathbf{W}^{k+1} = \underset{\mathbf{W}}{\operatorname{argmin}} \langle \mathbf{v}^k, \mathbf{W} - \mathbf{W}^k \rangle + \frac{1}{\alpha_k} \mathcal{D}_\phi(\mathbf{W}, \mathbf{W}^k) + \lambda_1 \|\mathbf{W}\|_1, \quad (14)$$

where $\widetilde{\nabla}_k$ is mini-batch gradient of $F(\mathbf{W})$. Omitting all the constants, the subproblem takes the form of:

$$\mathbf{W}^{k+1} = \underset{\mathbf{W}}{\operatorname{argmin}} \phi(\mathbf{W}) + \langle \mathbf{p}^k, \mathbf{W} \rangle + \alpha_k \lambda_1 \|\mathbf{W}\|_1, \quad (15)$$

where $\mathbf{p}^k = \alpha_k \mathbf{v}^k - \nabla \phi(\mathbf{W}^k)$. Here we adopt the kernel function $\phi(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|^2 + \frac{\delta}{r} \|\mathbf{W}\|^r$ ($r \geq 2$) for training neural networks, and then we have an explicit solution for (15) in Proposition 27.

Proposition 27 *Given $\mathbf{p}^k \in \mathbb{R}^d$, positive constant α_k, λ , and the kernel function $\phi(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|^2 + \frac{\delta}{r} \|\mathbf{W}\|^r$ ($r \geq 2, \delta > 0$). The solution of the subproblem (15) is given by*

$$\mathbf{W}^{k+1} = -t^* \mathbf{p}^+,$$

where t^* is the unique positive real root of the equation

$$(\delta \|\mathbf{p}^+\|^{r-2}) t^{r-1} + t - 1 = 0, \quad (16)$$

and \mathbf{p}^+ is given by

$$\mathbf{p}^+ = \underset{\mathbf{p}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{p} - \mathbf{p}^k\|^2 + \alpha_k \lambda \|\mathbf{p}\|_1 \right\}$$

which has an explicit expression given by $\mathbf{p}_j^+ = \operatorname{sign}(\mathbf{p}_j^k) \max(|\mathbf{p}_j^k| - \alpha_k \lambda, 0)$ for the j -th coordinate.

Proof The optimality condition of (15) is given by

$$0 = \mathbf{W}^{k+1} (1 + \delta \|\mathbf{W}^{k+1}\|^{r-2}) + \mathbf{p}^k + \alpha_k \lambda \mathbf{\Gamma}^k, \text{ where } \mathbf{\Gamma}^k \in \partial \|\cdot\|_1(\mathbf{W}^{k+1}).$$

Let $\mathbf{p}^+ = \mathbf{p}^k + \alpha_k \lambda \mathbf{\Gamma}^k$. By the optimality condition, we have $\mathbf{W}^{k+1} = -t \mathbf{p}^+$ for some positive scalar t , and

$$(-t - \delta \|\mathbf{p}^+\|^{r-2} t^{r-1} + 1) \mathbf{p}^+ = 0.$$

If $\mathbf{p}^+ \neq 0$, then $\delta \|\mathbf{p}^+\|^{r-2} t^{r-1} + t - 1 = 0$. If $\mathbf{p}^+ = 0$, then $\mathbf{W}^{k+1} = -t \mathbf{p}^+ = 0$. Since $t > 0$, then we have $\partial \|\cdot\|_1(\mathbf{W}^{k+1}) = \partial \|\cdot\|_1(-t \mathbf{p}^+) = -\partial \|\cdot\|_1(\mathbf{p}^+)$. Recall the definition of \mathbf{p}^+ , we have

$$\mathbf{p}^+ = \mathbf{p}^k + \alpha_k \lambda \mathbf{\Gamma}^k \in \mathbf{p}^k - \alpha_k \lambda \partial \|\cdot\|_1(\mathbf{p}^+),$$

which is sufficient and necessary optimality condition of the convex optimization problem:

$$\mathbf{p}^+ = \underset{\mathbf{p}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{p} - \mathbf{p}^k\|^2 + \alpha_k \lambda \|\mathbf{p}\|_1 \right\}.$$

This completes the proof by noting the the above minimization problem is the well-known soft threshold operator, see for example (Friedman et al., 2010). \blacksquare

In the absence of L_1 -regularization, that is, when $\lambda_1 = 0$, then $\mathbf{p}^+ = \mathbf{p}^k$ and the update formula for MSBPG at the k -th iteration simplifies to $\mathbf{W}^{k+1} = -t^* \mathbf{p}^k$, where t^* is the positive root of the equation (16). In this case, $\mathbf{W}^{k+1} = t^*(\nabla \phi(\mathbf{W}^k) - \alpha_k \mathbf{v}^k)$.

Furthermore, if we choose the kernel function simply as the square of Euclidean distance, i.e. $\delta = 0$, then SBPG reduces to SGD with momentum. Specifically, we have $t^* = 1$ and the update $\mathbf{W}^{k+1} = \mathbf{W}^k - \alpha_k \mathbf{v}^k$.

Determining degree of kernel function We now turn our attention to selecting an appropriate parameter r for the kernel function. Intuitively, in order to bound the Hessian of the loss function in (11), particularly when the number of layers L in (12) is large, r should also be chosen larger, so that $\nabla^2 F \preceq \frac{1}{\alpha} \nabla^2 \phi$ holds globally for some $\alpha > 0$. However, this can lead to numerical issues when computing $\|\mathbf{W}\|^{r-2}$. This problem can be avoided if the deep neural network exhibits some special structure such that a moderate r can make $F(\mathbf{W})$ smooth adaptable with respect to $\phi(\mathbf{W})$. For simplicity of analysis, we assume all the given labels y_i as zero and consider a sum of squares loss function. Then, we have a two-layer model defined as follows:

$$\min_{\mathbf{W}=(\mathbf{u},\mathbf{v})} F(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^N \left(\|\sigma(\text{Mat}(\mathbf{u})(g_i(\mathbf{v})))\|^2 \right) + \frac{\lambda_2}{2} (\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2), \quad (17)$$

where $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^{km}$, $g_i : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, $\text{Mat}(\mathbf{u}) \in \mathbb{R}^{k \times m}$ and $\sigma(\cdot)$ is a coordinate-wise operator. Notably, any deep neural network can be reformulated as the two-layer model in (17). For instance, if we define $\mathbf{v} = (\mathbf{W}_1, \dots, \mathbf{W}_{L-1})$, $\mathbf{u} = \text{Vec}(\mathbf{W}_L)$, $g_i(\mathbf{W}_1, \dots, \mathbf{W}_{L-1}) = \sigma_{L-1}(\mathbf{W}_{L-1}(\dots(\sigma_1(\mathbf{W}_1 \mathbf{x}_i))\dots))$, then model (12) can be transformed into (17). We make the following assumptions in this section, which guarantees that we can find a polynomial kernel function ϕ with a moderate degree, such that F in (17) is smooth adaptable to ϕ .

Assumption 4 σ is twice differentiable and σ' and $\sigma \cdot \sigma''$ are globally bounded.

Assumption 5 Each g_i is twice differentiable. All partial derivatives of order zero, one, and two of g_i are globally bounded.

Remark 28 Now we give some remarks on the above assumptions.

1. Assumption 4 is typically valid for various commonly used smooth activation functions. For example, the sigmoid function $\sigma(x) = \frac{1}{1+e^{-x}}$ satisfies global boundedness for both σ and σ'' . Certain activation function may not have bounded function value, such as GELU (Hendrycks and Gimpel, 2016), which takes the formulation of $\sigma(x) = x\Phi(x)$ where Φ is the standard Gaussian cumulative distribution function. Nonetheless, the product $\sigma \cdot \sigma''$ is globally bounded. Another type of activation function satisfying Assumption 4 is the smoothed ReLU function, for example, the following smoothed ReLU function, which will be considered in our numerical experiments:

$$\sigma_\epsilon(x) = \begin{cases} 0 & x \leq 0 \\ x^3 \left(\frac{1}{\epsilon^2} - \frac{x}{2\epsilon^3} \right) & 0 < x \leq \epsilon \\ x - \frac{\epsilon}{2} & x > \epsilon. \end{cases}$$

We observe that as ϵ tends to zero, σ_ϵ converges to the ReLU function. It is straightforward to verify that $\sigma_\epsilon \cdot \sigma''_\epsilon$ is globally bounded. Specifically, $\frac{3}{4}$ is a uniform bound on $\sigma_\epsilon \cdot \sigma''_\epsilon$ for $\epsilon \in (0, \frac{1}{2})$.

2. In many popular neural network frameworks, batch normalization (BN) layers (Ioffe and Szegedy, 2015) are often used before the fully connected layers. For example, in the VGG (Simonyan and Zisserman, 2014) and ResNet (He et al., 2016), BN layers are usually used before the last linear layer. In this case, we can treat all layers except the last one as one layer, which can be modeled as (17). It is expected that the BN layer can make the function g_i sufficiently smooth, thereby satisfying Assumption 5.

By applying the chain rule, we can compute the Hessian of F and determine a suitable degree parameter r in the kernel function, which will ensure that $\nabla^2 F$ is bounded by $\nabla^2 \phi$ globally. Consequently, F is smooth adaptable with respect to ϕ . In order to compute the Hessian of F , two formulas are required, which can be verified directly.

Lemma 29 *Let $\mathbf{u} \in \mathbb{R}^{km}$, $\mathbf{g} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^k$. Consider two linear maps: $\mathbf{u} \mapsto \text{Mat}(\mathbf{u})\mathbf{g}$ and $\mathbf{u} \mapsto \mathbf{A}(\text{Mat}(\mathbf{u}))^T \mathbf{b}$, then, the Jacobian of the two maps are given by*

$$\begin{aligned} J_{\mathbf{u}} [\text{Mat}(\mathbf{u})\mathbf{g}] &= \mathbf{g}^T \otimes \mathbb{I}_k, \\ J_{\mathbf{u}} [\mathbf{A}(\text{Mat}(\mathbf{u}))^T \mathbf{b}] &= \mathbf{A} \otimes \mathbf{b}^T. \end{aligned}$$

Proposition 30 *Suppose Assumptions 4 and 5 hold. Then, for any given $\delta > 0$ and any $r \geq 4$, the function F defined in (17) is smooth adaptable with respect to $\phi(\mathbf{W}) = \frac{1}{2}\|\mathbf{W}\|^2 + \frac{\delta}{r}\|\mathbf{W}\|^r$.*

Proof We denote $\text{Mat}(\mathbf{u})$ by \mathbf{M} . The Jacobian of g is denoted by Jg , while its transpose is denoted by $J^T g$. \mathbb{I}_k is $k \times k$ identity matrix. We only need to prove the single sample case, i.e. $N = 1$. Using Lemma 29, we can compute the Jacobian and Hessian of F as follows:

Jacobian of F :

$$\begin{aligned} \frac{\partial F}{\partial \mathbf{u}} &= (g(\mathbf{v}) \otimes \mathbb{I}_k) [\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma(\mathbf{M}g(\mathbf{v}))] + \lambda_2 \mathbf{u}, \\ \frac{\partial F}{\partial \mathbf{v}} &= J^T g(\mathbf{v}) \mathbf{M}^T [\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma(\mathbf{M}g(\mathbf{v}))] + \lambda_2 \mathbf{v}. \end{aligned} \tag{18}$$

Hessian of F :

$$\begin{aligned} \frac{\partial^2 F}{\partial \mathbf{u}^2} &= (1) + (2) + \lambda_2 \mathbb{I}_{km}, \\ \text{where (1)} &= (g(\mathbf{v}) \otimes \mathbb{I}_k) \text{Diag}(\sigma(\mathbf{M}g(\mathbf{v})) \circ \sigma''(\mathbf{M}g(\mathbf{v}))) (g^T(\mathbf{v}) \otimes \mathbb{I}_k) \\ (2) &= (g(\mathbf{v}) \otimes \mathbb{I}_k) \text{Diag}(\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma'(\mathbf{M}g(\mathbf{v}))) (g^T(\mathbf{v}) \otimes \mathbb{I}_k). \end{aligned} \tag{19}$$

$$\begin{aligned} \frac{\partial^2 F}{\partial \mathbf{u} \partial \mathbf{v}} &= (1) + (2) + (3), \\ \text{where (1)} &= (J^T g(\mathbf{v})) \otimes [\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma'(\mathbf{M}g(\mathbf{v}))]^T \\ (2) &= J^T g(\mathbf{v}) \mathbf{M}^T \text{Diag}[\sigma(\mathbf{M}g(\mathbf{v})) \circ \sigma''(\mathbf{M}g(\mathbf{v}))] (g^T(\mathbf{v}) \otimes \mathbb{I}_k) \\ (3) &= J^T g(\mathbf{v}) \mathbf{M}^T \text{Diag}[\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma'(\mathbf{M}g(\mathbf{v}))] (g^T(\mathbf{v}) \otimes \mathbb{I}_k). \end{aligned} \tag{20}$$

$$\frac{\partial^2 F}{\partial \mathbf{v}^2} = (1) + (2) + (3) + \lambda_2 \mathbb{I}_n,$$

$$\text{where (1)} = D^2 g(\mathbf{v}) [\mathbf{M}^T [\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma(\mathbf{M}g(\mathbf{v}))]] = \sum d_i \nabla^2 g_i(\mathbf{v}) \quad (21)$$

$$(2) = J^T g(\mathbf{v}) \mathbf{M}^T \text{Diag}[\sigma(\mathbf{M}g(\mathbf{v})) \circ \sigma''(\mathbf{M}g(\mathbf{v}))] \mathbf{M} J g(\mathbf{v})$$

$$(3) = J^T g(\mathbf{v}) \mathbf{M}^T \text{Diag}[\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma'(\mathbf{M}g(\mathbf{v}))] \mathbf{M} J g(\mathbf{v}),$$

where $\mathbf{d} = \mathbf{M}^T [\sigma'(\mathbf{M}g(\mathbf{v})) \circ \sigma(\mathbf{M}g(\mathbf{v}))]$. Now, we are ready to prove this proposition. For any $\mathbf{w} \in \mathbb{R}^{km+n}$ and $\mathbf{h} = [\mathbf{h}^u; \mathbf{h}^v] \in \mathbb{R}^{km+n}$, it suffices to prove that $\langle \nabla^2 F(\mathbf{w}) \mathbf{h}, \mathbf{h} \rangle = \mathcal{O}(\langle \nabla^2 \phi(\mathbf{w}) \mathbf{h}, \mathbf{h} \rangle)$. From (19)(20)(21) and Assumption 4, 5, we can easily get $\langle \nabla^2 F(\mathbf{w}) \mathbf{h}, \mathbf{h} \rangle = \mathcal{O}((1 + \|\mathbf{w}\|^2) \|\mathbf{h}\|^2)$. On the other hand, $\nabla^2 \phi(\mathbf{w}) = I(1 + \|\mathbf{w}\|^{r-2}) + (r-2) \|\mathbf{w}\|^{r-4} \mathbf{w} \mathbf{w}^T$. Hence $\langle \nabla^2 \phi(\mathbf{w}) \mathbf{h}, \mathbf{h} \rangle \geq (1 + \|\mathbf{w}\|^{r-2}) \|\mathbf{h}\|^2$. So, we only require $r-2 \geq 2$. This completes the proof. \blacksquare

Layerwise kernel function In Proposition 27, we use the kernel function $\phi(\mathbf{W}) = \frac{1}{2} \|\mathbf{W}\|^2 + \frac{\delta}{r} \|\mathbf{W}\|^r$, which applies the same Bregman distance across all layers of the deep neural network. However, different layers exhibit distinct geometric properties (You et al., 2019), and computing $\|\mathbf{W}\|^r$ with $r > 2$ may result in numerical instability for neural networks with millions of parameters, such as in VGG (Simonyan and Zisserman, 2014). To take advantage of the layerwise structure of neural networks, we design a layerwise kernel function for a L -layer neural network as follows:

$$\phi(\mathbf{W}) = \sum_{i=1}^L \phi_i(\mathbf{W}_i), \quad \phi_i(\mathbf{W}_i) = \frac{1}{2} \|\mathbf{W}_i\|^2 + \frac{\delta}{r} \|\mathbf{W}_i\|^r. \quad (22)$$

Note that δ and r can vary from layer to layer, here we take the same δ and r for all layers for simplicity. With this structure, the Bregman distance takes the form $\mathcal{D}_\phi = \sum_{i=1}^L \mathcal{D}_{\phi_i}$. By incorporating this layerwise Bregman distance into the subproblem (14), our MSBPG algorithm can be implemented in a layerwise manner. Details of the implementation are provided in Algorithm 1.

Mitigating gradient explosion In the training of deep neural networks, gradient explosion is a common undesired phenomenon, where the gradients of the loss function grow exponentially from layer to layer, leading to numerical instability or even collapse of the training process (Hochreiter, 1991; Manchev and Spratling, 2020). The reasons for gradient explosion include selecting a large stepsize and choosing an improper initialization for the model's parameters (Pascanu et al., 2013). In the following, we will show that MSBPG provides a novel approach to mitigate gradient explosion. Considering MSBPG without L_1 -regularization, the update rule is given by:

$$\mathbf{W}_i^{k+1} = -t_i^k \mathbf{p}_i^k = t_i^k \left((1 + \delta \|\mathbf{W}_i^k\|^{r-2}) \mathbf{W}_i^k - \alpha_k \mathbf{v}_i^k \right), \quad (23)$$

where $t_i^k \in (0, 1)$ is the unique positive root of

$$\left(\delta \left\| (1 + \delta \|\mathbf{W}_i^k\|^{r-2}) \mathbf{W}_i^k - \alpha_k \mathbf{v}_i^k \right\|^{r-2} \right) t^{r-1} + t - 1 = 0. \quad (24)$$

Algorithm 1 Momentum based Stochastic Bregman Proximal Gradient (MSBPG) for training neural networks

- 1: **Input:** Total number of iterations K , stepsize $\alpha_k > 0$, momentum parameter $\beta_k \in (0, 1)$, $\delta > 0$ and integer $r \geq 4$ to determine the kernel function ϕ , and λ_1 .
 - 2: **Initialize:** Set $\mathbf{W} = \mathbf{W}^0$, $\mathbf{v}^0 = 0$.
 - 3: **for** $k = 0, \dots, K - 1$ **do**
 - 4: Compute mini-batch gradient $\tilde{\nabla}_k$ of F ;
 - 5: Compute SMAE: $\mathbf{v}^k = (1 - \beta_k)\mathbf{v}^{k-1} + \beta_k\tilde{\nabla}_k$;
 - 6: **for** $i = 1, \dots, L$ **do**
 - 7: $\mathbf{p}_i^k = \alpha_k \mathbf{v}_i^k - \nabla \phi(\mathbf{W}_i^k)$;
 - 8: $\mathbf{p}_i^+ = \operatorname{argmin}_{\mathbf{p}_i} \{\frac{1}{2}\|\mathbf{p}_i - \mathbf{p}_i^k\|^2 + \alpha_k \lambda_1 \|\mathbf{p}_i\|_1\}$;
 - 9: Solve $(\delta \|\mathbf{p}_i^+\|^{r-2})t_i^{r-1} + t_i - 1 = 0$ to get t_i^k ;
 - 10: $\mathbf{W}_i^{k+1} = -t_i^k \mathbf{p}_i^+$;
 - 11: **end for**
 - 12: **end for**
 - 13: **Output:** $\mathbf{W}^1, \dots, \mathbf{W}^K$.
-

Combining (23) and (24), we have the following equivalent implicit update scheme for the i -th layer:

$$\mathbf{W}_i^{k+1} = \frac{1 + \delta \|\mathbf{W}_i^k\|^{r-2}}{1 + \delta \|\mathbf{W}_i^{k+1}\|^{r-2}} \mathbf{W}_i^k - \frac{\alpha_k}{1 + \delta \|\mathbf{W}_i^{k+1}\|^{r-2}} \mathbf{v}_i^k. \quad (25)$$

It is observed in practice that with large stepsize or large initial point, the gradient \mathbf{v}_i^k tends to explode if no scaling or clipping is applied, while the norm of the weight $\|\mathbf{W}_i^{k+1}\|$ also tends to be large. In (25), scaling the gradient by $\frac{1}{1 + \delta \|\mathbf{W}_i^{k+1}\|^{r-2}}$ prevents the weight \mathbf{W}_i^{k+1} from moving too drastically in the direction of the gradient, thereby controlling the rapid growth of its norm. At the same time, if the norm $\|\mathbf{W}_i^{k+1}\|$ does not change significantly, the coefficient of \mathbf{W}_i^k in (25) will remain approximately 1. Thus, the implicit update (25) provides automatic scaling of the gradient, effectively mitigating the rapid growth of the weight and preventing subsequent gradient explosion.

Experimental results in Section 5.2 indeed verify MSBPG’s ability to mitigate gradient explosion for training deep neural networks. An intuitive illustration of SBPG’s “pull-back” ability is given in Figure 12 in Appendix D, and this “pull-back” ability originates from the Bregman proximity model and the polynomial kernel function we adopt.

5. Numerical experiments

In this section, we conduct numerical experiments to demonstrate the effectiveness and robustness of MSBPG in comparison to some commonly used optimizers in deep learning. We assess the impact of stepsize and initial point selection on the performance of our method. Our experiments consist of two parts. In the first part, we use a quadratic inverse problem as a toy example to illustrate the capabilities of vanilla SBPG. The second part is the main focus of this section, where we evaluate the performance of MSBPG in training deep neural networks. The experiments for the quadratic inverse problem are conducted

using MATLAB R2021b on a Windows workstation equipped with a 12-core Intel Xeon E5-2680 @ 2.50GHz processor and 128GB of RAM. For the deep learning experiments, we conducted the experiments using PyTorch running on a single RTX3090 GPU.

5.1 Quadratic inverse problem

The quadratic inverse problem, as formulated in (Bolte et al., 2018), is given by:

$$\min \left\{ \Phi(\mathbf{x}) := \frac{1}{4} \underbrace{\sum_{i=1}^n (\langle A_i \mathbf{x}, \mathbf{x} \rangle - b_i)^2}_{F(\mathbf{x})} + \lambda R(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d \right\},$$

which has practical applications (Beck and Eldar, 2013), including the phase retrieval problem as a special case (Luke, 2017). In this experiment, we consider the L_1 regularization $R(\mathbf{x}) = \|\mathbf{x}\|_1$ with $\lambda = 1 \times 10^{-3}$, and solve the quadratic inverse problem using SBPG and stochastic (Euclidean) proximal gradient (SPG) method (Bertsekas, 2011). Notably, SPG is a special case of SBPG, where $\phi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$. Since the smooth term in the objective function $F(\mathbf{x})$ does not admit a globally Lipschitz continuous gradient, we employ the kernel function $\phi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2 + \frac{1}{r}\|\mathbf{x}\|^r$ with $r = 4$. It has been shown in (Lu et al., 2018) that any $r \geq 4$ guarantees that F is ϕ -smooth adaptable globally. Moreover, according to (Bolte et al., 2018), the smooth adaptable constant L_F can be chosen such that $L_F \geq \sum_{i=1}^n (3\|A_i\|^2 + \|A_i\| |b_i|)$ for $r = 4$. In this experiment, we randomly generate the data by the following MATLAB commands:

```
ai = randn(d, 1); Ai = ai * ai';
x_true = sprandn(d, 1, density_x); b_i = x_true' * (Ai * x_true);
```

The true solution for the quadratic inverse problem is chosen as a sparse vector \mathbf{x}^* that satisfies $\langle A_i \mathbf{x}^*, \mathbf{x}^* \rangle = b_i$ for $i = 1, \dots, n$. We set the mini-batch size for all algorithms to be $m = 1$. To evaluate the effectiveness of each algorithm, we use the following criterion that takes into account the possibility of stationary points being local minimum or saddle points:

$$\epsilon_k = \max \left\{ \|\mathcal{G}(\mathbf{x}^k)\|, \epsilon_{\text{obj}} := \frac{\text{obj}_k - \text{obj}_*}{1 + \text{obj}_*} \right\}, \quad (26)$$

where $\text{obj}_k = \Phi(\mathbf{x}^k)$ and $\text{obj}_* = \Phi(\mathbf{x}^*)$. The term $\|\mathcal{G}(\mathbf{x}^k)\|$ measures the stationarity of the solution, while a small ϵ_{obj} indicates that the solution is a "nearly" global minimum.

We conduct experiments on a problem with data size $d = 100$ and $\text{density}_x = 0.05$. All methods are run until they reach an accuracy of $\epsilon_k \leq 0.01$ within a time limit of 30 seconds. To ensure statistical significance, we run each algorithm 10 times and report the median value. The results are presented in Figure 2. For Figures 2(a), we randomly select initial points within a ball centered at the origin with radius 1×10^{-2} . We use the stepsize schedule of $\alpha_k = \max \left\{ 10^{-4}, \frac{\alpha_0}{\sqrt{1+k}} \right\}$, where α_0 is the initial stepsize. For Figure 2(b), we set constant stepsize schedule 1×10^{-3} . For Figures 2(c), we randomly select initial points within a ball centered at the origin with radius 1×10^{-2} . We use a constant stepsize

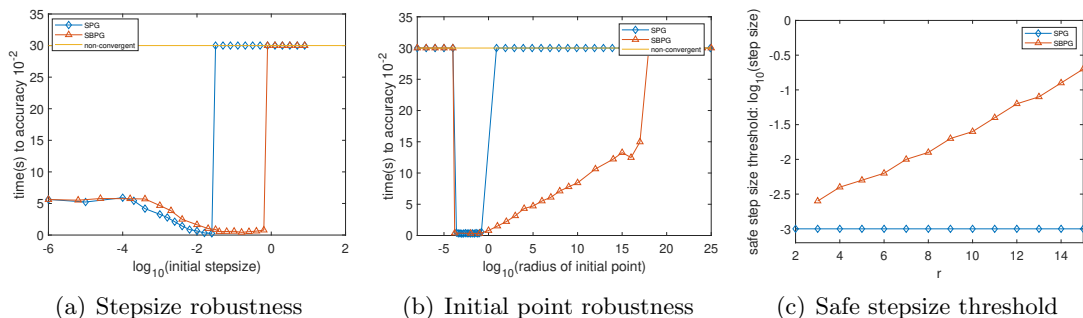


Figure 2: Comparison of SBPG and SPG in terms of their robustness with respect to stepsize and initial point selection. A method is considered non-convergent if it fails to reach an accuracy of $\epsilon_k < 10^{-2}$ within 30 seconds or if it collapses. Generally, choosing large stepsize and large radius for the initial point can cause an algorithm to collapse. The safe stepsize threshold is the maximum stepsize (constant schedule) that a method does not collapse. We run 10 tests for each algorithm and report the median of the results.

schedule. To prevent excessively small stepsizes that can slow down all methods, we set a lower bound for the stepsize.

Figure 2(a) demonstrates that SBPG has a wider range of convergent stepsizes than SPG, indicating that SBPG is more robust in terms of stepsize selection. The effect of the initial stepsize on the performance of the algorithms is also depicted in this figure. Figure 2(b) highlights that SBPG is significantly more robust than SPG with respect to initial point selection, showing high resilience and preventing the training process from collapsing. Additionally, Figure 2(c) illustrates that increasing the degree r in the kernel function raises the threshold for safe stepsizes. These observations are explained in Section 4. When a large stepsize or a large initial point radius leads to a potential gradient explosion, the Bregman proximal mapping effectively pulls back the iterate, guiding it toward a more stable region and preventing gradient explosion.

5.2 Deep neural network

To evaluate MSBPG’s performance in training deep neural networks, we consider a model with L_2 regularization to enhance generalization:

$$\min_{\mathbf{W}} \underbrace{\frac{1}{N} \sum_{i=1}^N \mathcal{L}(\mathcal{DNN}(\mathbf{W}, \mathbf{x}_i), y_i)}_{F(\mathbf{W})} + \lambda_2 \|\mathbf{W}\|_2^2. \quad (27)$$

We employ MSBPG to solve this large-scale problem. Following AdamW (Loshchilov and Hutter, 2017), we employ the strategy of decoupled weight decay. We use MSBPG to solve this large-scale optimization problem. Following the AdamW approach (Loshchilov and Hutter, 2017), we implement a decoupled weight decay strategy. Specifically, we compute

the stochastic gradient only for the loss function F , treating F and L_2 regularization separately. After performing the Bregman proximal mapping based on the stochastic gradient, we apply a weight decay step.

The detailed algorithm is summarized in Algorithm 2. At iteration k , MSBPG first uses automatic differentiation to compute the mini-batch gradient $\tilde{\nabla}_k$ of F . Then, it maintains a bias-corrected gradient estimator $\bar{\mathbf{v}}^k$ (Kingma and Ba, 2014) and uses it to calculate the layerwise \mathbf{p}_i^k . With \mathbf{p}_i^k , MSBPG solves a univariate equation to get t_i^k and updates the weight of the i -th layer to \mathbf{W}_i^k . In the end, MSBPG conducts decoupled weight decay.

Algorithm 2 MSBPG with L_2 regularization

- 1: **Input:** Total number of training epochs K , momentum coefficient β , stepsize α_k , weight decay coefficient λ_2 , δ and r to determine the kernel function ϕ .
 - 2: **Initialize:** Set $\mathbf{W} = \mathbf{W}^0$, $\mathbf{v}^0 = \mathbf{0}$.
 - 3: **for** $k = 1, \dots, K$ **do**
 - 4: Compute mini-batch gradient $\tilde{\nabla}_k$;
 - 5: $\mathbf{v}^k = \beta \mathbf{v}^{k-1} + (1 - \beta) \tilde{\nabla}_k$, $\bar{\mathbf{v}}^k = \mathbf{v}^k / (1 - \beta^k)$;
 - 6: **for** $i = 1, \dots, L$ **do**
 - 7: $\mathbf{p}_i^k = \alpha_k \bar{\mathbf{v}}_i^k - \nabla \phi(\mathbf{W}_i^{k-1})$;
 - 8: Solve $(\delta \|\mathbf{p}_i^k\|^{r-2}) t_i^{r-1} + t_i - 1 = 0$ to get t_i^k ;
 - 9: $\tilde{\mathbf{W}}_i^k = -t_i^k \mathbf{p}_i^k$;
 - 10: **end for**
 - 11: $\mathbf{W}^k = \tilde{\mathbf{W}}^k - \alpha_k \lambda_2 \mathbf{W}^{k-1}$;
 - 12: **end for**
 - 13: **Output:** $\mathbf{W}^1, \dots, \mathbf{W}^K$
-

We conducted experiments on several representative benchmarks, including VGG16 (Simonyan and Zisserman, 2014), ResNet34 (He et al., 2016) on CIFAR10 dataset (Krizhevsky et al., 2009), ResNet34 (He et al., 2016), DenseNet121 (Huang et al., 2017) on CIFAR100 dataset (Krizhevsky et al., 2009), and LSTMs (Hochreiter and Schmidhuber, 1997) on the Penn Treebank dataset (Marcinkiewicz, 1994). We compare MSBPG with the most popular optimization algorithms used for training neural networks, including SGD (Sutskever et al., 2013), Adam (Kingma and Ba, 2014), and AdamW (Loshchilov and Hutter, 2017). Experimental results show that MSBPG has good convergence performance and generalization capacity for both the task that SGD dominates (image classification with CNNs) and the task that Adam dominates (language modeling with LSTMs). We also conducted experiments to compare MSBPG with SGD on different initial stepsizes and different scales of the initial point. Our experimental results demonstrate the robustness of MSBPG in training neural networks.

To further evaluate the performance of MSBPG on more recent neural network architectures, we conducted additional experiments, with the results presented in Appendix C.

Before getting into the details of our experiments, we first clarify the activation function. The commonly used ReLU activation function in VGG, ResNet, and DenseNet is defined as $\text{ReLU}(x) = \max(0, x)$, which is not continuously differentiable. To address this, we design a smooth approximation of ReLU with a parameter ϵ , which is twice continuously

differentiable and satisfies our Assumption 4:

$$\sigma_\epsilon(x) = \begin{cases} 0 & x \leq 0 \\ x^3 \left(\frac{1}{\epsilon^2} - \frac{x}{2\epsilon^3} \right) & 0 < x \leq \epsilon \\ x - \frac{\epsilon}{2} & x > \epsilon. \end{cases}$$

The gradient of this activation function is given by:

$$\sigma'_\epsilon(x) = \begin{cases} 0 & x \leq 0 \\ x^2 \left(\frac{3}{\epsilon^2} - \frac{2x}{\epsilon^3} \right) & 0 < x \leq \epsilon \\ 1 & x > \epsilon. \end{cases}$$

As ϵ tends to 0, this smooth activation function converges to the standard ReLU function. We conducted experiments using VGG16 on the CIFAR-10 dataset, replacing all activation functions in VGG16 with σ_ϵ . As shown in Figure 4, our algorithm MSBPG’s performance does not degrade as ϵ tends to 0. Therefore, in the subsequent experiments, we use the original neural network architectures with the ReLU activation function to evaluate our method MSBPG.

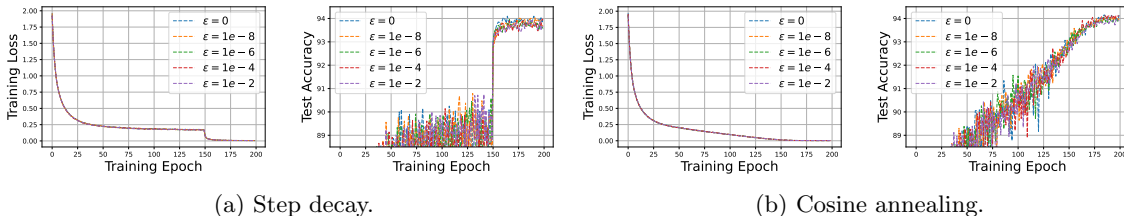


Figure 3: Training loss and test accuracy (%) of VGG16 on CIFAR10 dataset under two frequently used training settings. Here the activation function of VGG16 adopts smoothed ReLU activation function σ_ϵ with different choices of ϵ ($\epsilon = 0$ denotes adopting the original ReLU activation function).

CNNs on image classification We conducted experiments with VGG16 and ResNet34 on the CIFAR-10 dataset, and ResNet34 and DenseNet121 on the CIFAR-100 dataset. SGD usually has better generalization performance than adaptive gradient algorithms such as Adam and AdamW when training CNNs on image classification tasks. For our experiments, we utilized two common training strategies: reducing the stepsize to 10% of its original value near the end of training (Zhuang et al., 2020; Chen et al., 2021; Luo et al., 2019), and using a cosine annealing schedule for stepsizes (Loshchilov and Hutter, 2016, 2017). These strategies are designed to accelerate convergence and allow for a fair comparison of the generalization capacities of different optimizers. We use the default training hyperparameters of SGD, Adam, and AdamW in these settings (He et al., 2016; Zhuang et al., 2020; Chen et al., 2021), and set MSBPG’s learning rate (initial stepsize) as 0.1, momentum coefficient β as 0.9, weight decay coefficient λ_2 as 1×10^{-3} . For the layerwise kernel function $\phi_i(\mathbf{W}_i) = \frac{1}{2} \|\mathbf{W}_i\|^2 + \frac{\delta}{r} \|\mathbf{W}_i\|^r$, we set $r = 4$, $\delta = 1 \times 10^{-2}$ for VGG16 and $r = 6$, $\delta = 1 \times 10^{-3}$ for

ResNet34 on CIFAR10 dataset, and $r = 4$, $\delta = 1 \times 10^{-2}$ for ResNet34 and $r = 4$, $\delta = 1 \times 10^{-3}$ for DenseNet121 on CIFAR100 dataset.

Our convergence analysis demonstrates that MSBPG converges to a stationary point. As shown by (Lee et al., 2019; Panageas et al., 2019), the Bregman gradient method almost always converges to a local minimum for loss functions with the strict saddle property. When the neural network is highly overparameterized and exhibits a benign nonconvexity in the search region, stochastic first-order methods tend to find the global minimum, where the loss function value is 0.

From the experimental results in Figure 4, 5, 6, 7, we can observe that MSBPG attains almost zero training loss in all the training settings. This implies that our method can find the global minimum in these instances. Furthermore, MSBPG consistently achieves the highest test accuracy for all experimental settings and attains at least 0.5% test accuracy improvement compared with the second-best optimization algorithm. This generalization advantage of MSBPG can be attributed to the Bregman proximity model we adopt. In Appendix B, we further discuss SBPG’s performance from the perspective of algorithmic stability, which can influence the generalization gap. We demonstrate that the high-order polynomial kernel used helps reduce the generalization gap bound in high-dimensional scenarios, which partially explains the good generalization performance of our proposed methods.

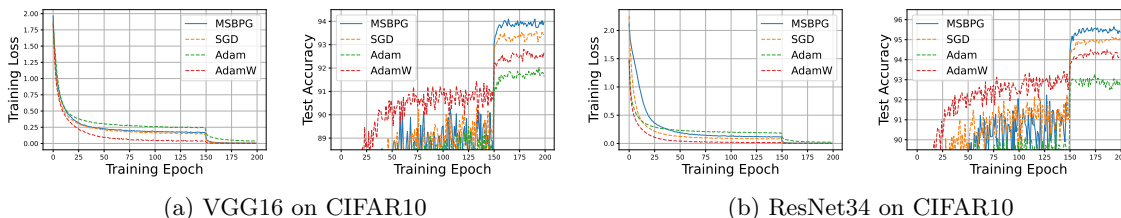


Figure 4: Training loss and test accuracy (%) of CNNs on CIFAR10 dataset with learning rate reduced to 0.1 times of the original value at the 150th epoch.

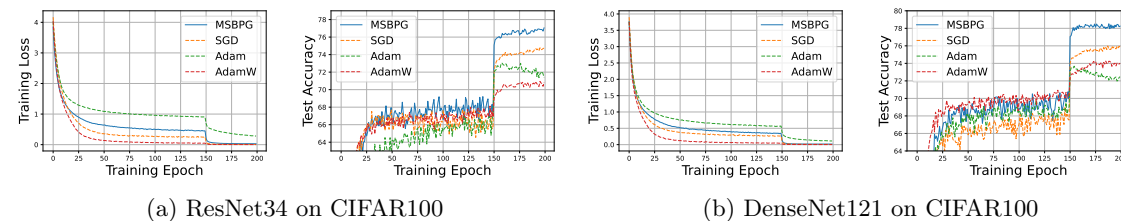


Figure 5: Training loss and test accuracy (%) of CNNs on CIFAR100 dataset with learning rate reduced to 0.1 times of the original value at the 150th epoch.

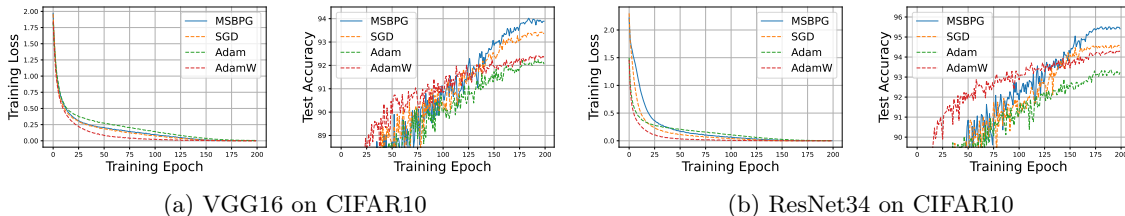


Figure 6: Training loss and test accuracy (%) of CNNs on CIFAR10 dataset with learning rate using the cosine annealing schedule.

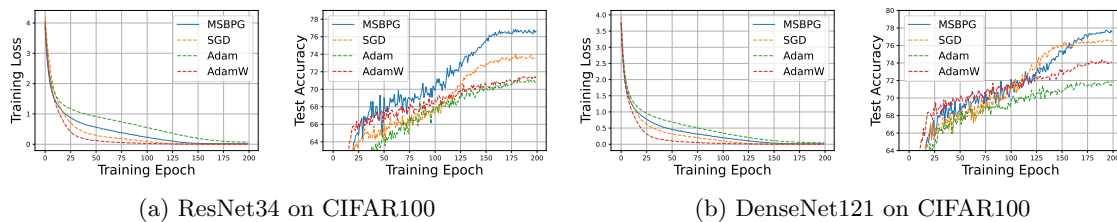


Figure 7: Training loss and test accuracy (%) of CNNs on CIFAR100 dataset with learning rate using the cosine annealing schedule.

LSTMs on language modeling To further evaluate the performance of MSBPG, we conducted experiments on LSTMs using the Penn Treebank dataset, reporting both training and test perplexity (lower is better). Adam is generally favored over SGD for language modeling tasks due to its better generalization capacity (Fu et al., 2016; Siami-Namini et al., 2019), and thus it is the default optimization algorithm for training LSTMs. We followed the standard experimental setup for training LSTMs (Zhuang et al., 2020; Chen et al., 2021), where the learning rate is reduced to 10% of its original value twice (at the 75th and 150th epochs). Additionally, we experimented with the cosine annealing learning rate schedule (Loshchilov and Hutter, 2016), which is commonly used in practice. For the training hyperparameters, we used the default settings for SGD, Adam, and AdamW when training 1-, 2-, and 3-layer LSTMs (Zhuang et al., 2020; Chen et al., 2021). For MSBPG, we set the learning rate to 25, 80, and 80 for 1-, 2-, and 3-layer LSTMs, respectively, with a momentum parameter $\beta = 0.9$, weight decay coefficient $\lambda_2 = 2 \times 10^{-6}$. For the layerwise kernel function $\phi_i(\mathbf{W}_i) = \frac{1}{2} \|\mathbf{W}_i\|^2 + \frac{\delta}{r} \|\mathbf{W}_i\|^r$, we set $r = 4$ and $\delta = 1 \times 10^{-6}$. From Figure 8 and Figure 9, we observe that MSBPG converges well on the training dataset for 1-, 2-, and 3-layer LSTMs across both training strategies. In contrast, SGD with the cosine annealing learning rate schedule fails to fully converge on the training dataset, as shown in Figure 9. Additionally, MSBPG consistently achieves lower test perplexity in all experiments, outperforming other methods by at least 1 unit. This excellent generalization capacity can be attributed to the Bregman proximity model employed by MSBPG. As an additional evaluation, we further assess the performance of MSBPG on a recently popular transformer

model, Transformer-XL (Dai et al., 2019), which is designed for long-sequence tasks. The results are provided in Appendix C.

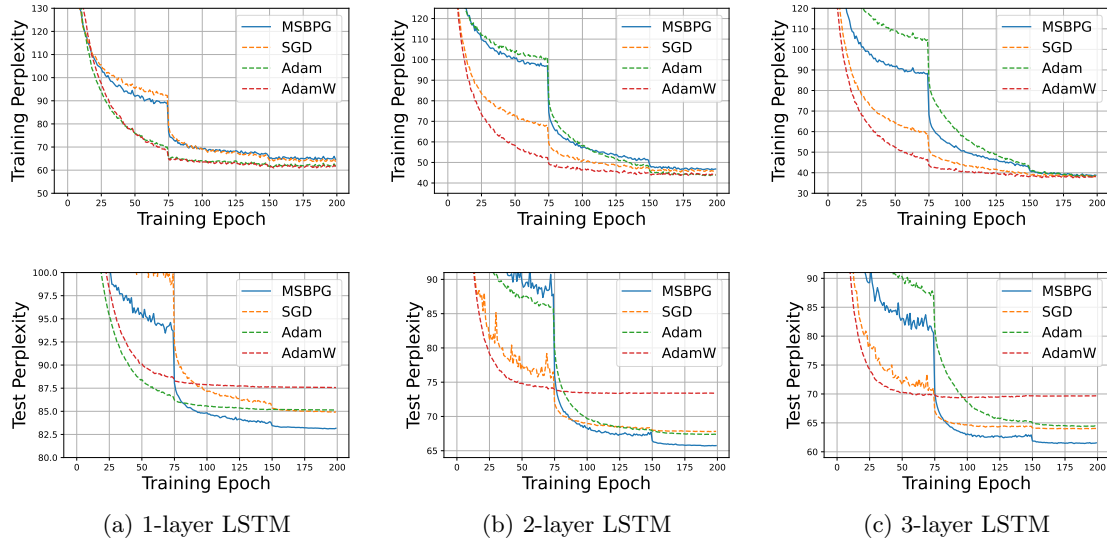


Figure 8: Training and test perplexity (lower is better) of LSTMs on Penn Treebank dataset with learning rate reduced to 0.1 times of the original value at the 75th epoch and 150th epoch.

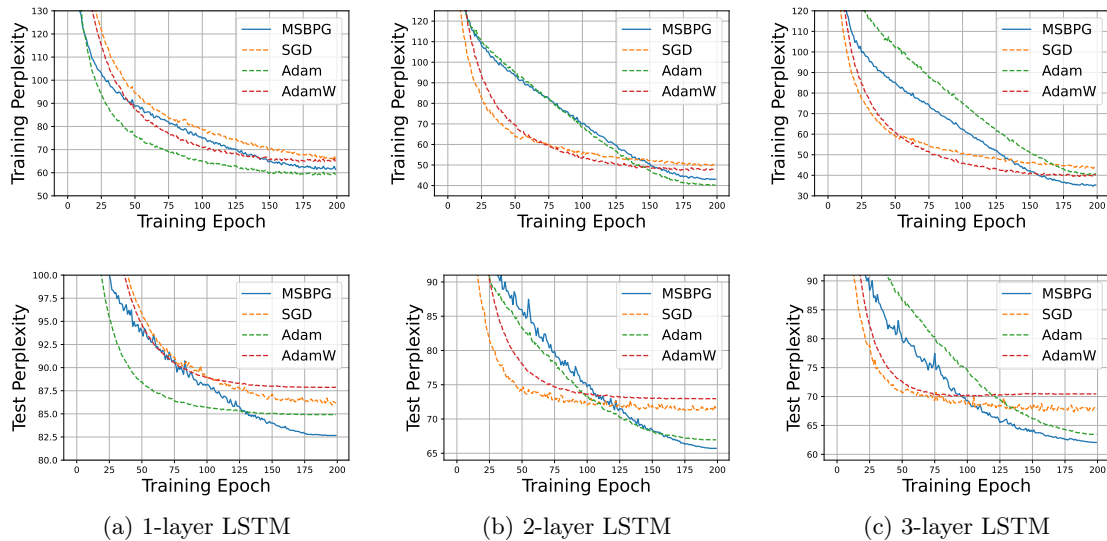


Figure 9: Training and test perplexity (lower is better) of LSTMs on Penn Treebank dataset with learning rate using the cosine annealing schedule.

Robustness to initial point scale and stepsize As demonstrated in Section 4, MSBPG can mitigate the issue of gradient explosion, which typically occurs when using large stepsizes or large initial point scales. To verify MSBPG’s robustness in training neural networks, we conducted experiments using VGG16 on the CIFAR-10 dataset. Specifically, we compared the performance of MSBPG and SGD under different initial point scales and stepsizes, as both algorithms share the same default learning rate of (1×10^{-1}) . Since adaptive gradient algorithms, such as Adam, use a different default learning rate scale (1×10^{-3}), they were not included in this comparison. For various initial point scales and stepsizes, we ran each optimization algorithm for 50 iterations and reported the best test accuracy. As shown in Figure 10, MSBPG is more robust than SGD to large initial points and stepsizes. Training deep neural networks, which have millions or billions of parameters, is highly sensitive to both the initial point scale and stepsize. From Figure 10, we can see that SGD fails to converge when the initial point scale is increased to 4.6 or when the stepsize is increased from 0.1 to 0.6. In contrast, MSBPG converges with an initial point scale as large as 20 and a stepsize as large as 5. This robustness of MSBPG can ease the tuning of hyperparameters for training neural networks, and can also make the training process more robust to noises and errors.

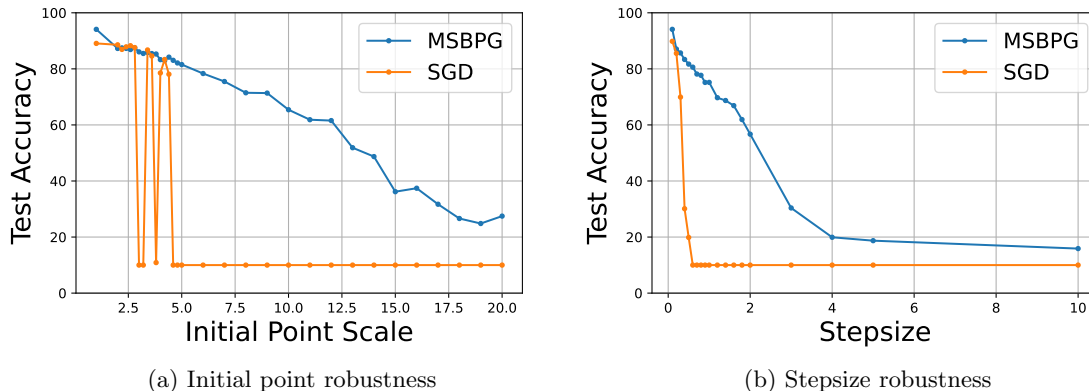


Figure 10: Test accuracy (%) of VGG16 on CIFAR10 dataset with different initial point scale and stepsize choice.

6. Conclusion

In this paper, we introduce a family of nonconvex stochastic Bregman proximal gradient (SBPG) methods to solve optimization problems without the Lipschitz smoothness assumption. By leveraging Bregman proximity measures, SBPG offers a more flexible and robust framework than classical stochastic gradient methods. We establish convergence results for the vanilla SBPG method in the nonconvex setting and propose a momentum-based variant, MSBPG, which improves convergence property by relaxing the mini-batch size requirement. Both methods achieve optimal sample complexity $\mathcal{O}(\epsilon^{-4})$, making them well-suited for large-scale problems. MSBPG is applied to training deep neural networks, where it mit-

igates gradient explosion and enhances generalization performance. Numerical experiments on sparse quadratic inverse problems and deep neural networks demonstrate MSBPG's robustness and superior performance compared to commonly used optimizers such as SGD, Adam, and AdamW. In conclusion, MSBPG provides an effective and efficient optimization approach for large-scale nonconvex problems, combining theoretical robustness with practical advantages. Future work can explore further refinements and applications in more complex machine learning tasks.

Appendix A. Proofs in Preliminaries

Proof of Lemma 8 First, we prove the uniqueness of the solution. Problem (4) is equivalent to the following problem:

$$\arg \min_{\mathbf{u} \in \bar{C}} \Psi(\mathbf{u}) := \alpha R(\mathbf{u}) + \phi(\mathbf{u}) - \langle \nabla \phi(\mathbf{x}), \mathbf{u} \rangle.$$

We have that

$$\Psi(\mathbf{u}) \geq \alpha R(\mathbf{u}) + \phi(\mathbf{u}) - \|\nabla \phi(\mathbf{x})\| \|\mathbf{u}\| \geq \|\mathbf{u}\| \left(\frac{\alpha R(\mathbf{u}) + \phi(\mathbf{u})}{\|\mathbf{u}\|} - \|\nabla \phi(\mathbf{x})\| \right).$$

As $\|\mathbf{u}\| \rightarrow \infty$, we have $\Psi(\mathbf{u}) \geq \|\mathbf{u}\| \left(\frac{\alpha R(\mathbf{u}) + \phi(\mathbf{u})}{\|\mathbf{u}\|} - \|\nabla \phi(\mathbf{x})\| \right) = \infty$, where we use the fact that ϕ is supercoercive and R is convex. Since Ψ is a proper lower-semicontinuous convex function, by the modern form of Weierstrass theorem (Rockafellar, 1997, Chapter 1), we know that the solution set of (4) is a nonempty compact set. Also note that Ψ is a strictly convex function, which implies the uniqueness of the solution. For any Legendre function ϕ , from (Rockafellar, 1997, Chapter 26), we have $\text{dom } \partial \phi = \text{int } \text{dom } \phi$ with $\partial \phi(\mathbf{x}) = \{\nabla \phi(\mathbf{x})\}$ for all $\mathbf{x} \in \text{int } \text{dom } \phi$. The optimality condition implies that $\partial \phi(\text{Prox}_{\alpha R}^{\phi}(\mathbf{x}))$ is nonempty, which automatically forces $\text{Prox}_{\alpha P}^{\phi}(\mathbf{x}) \in \text{int } \text{dom } \phi$. This completes the proof. \square

Proof of Proposition 12 Note that $\|\nabla \phi(\mathbf{x}^+) - \nabla \phi(\mathbf{x})\| \leq L_{\phi} \|\mathbf{x}^+ - \mathbf{x}\|$ and $\|\nabla F(\mathbf{x}^+) - \nabla F(\mathbf{x})\| \leq L_F L_{\phi} \|\mathbf{x}^+ - \mathbf{x}\|$. By the definition of \mathbf{x}^+ , we have

$$\nabla F(\mathbf{x}^+) - \nabla F(\mathbf{x}) + \frac{\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}^+)}{\alpha} \in \nabla F(\mathbf{x}^+) + \partial R(\mathbf{x}^+).$$

Thus, we obtain

$$\text{dist}(0, \partial \Phi(\mathbf{x}^+)) \leq \left\| \nabla F(\mathbf{x}^+) - \nabla F(\mathbf{x}) + \frac{\nabla \phi(\mathbf{x}) - \nabla \phi(\mathbf{x}^+)}{\alpha} \right\| \leq \left(L_F L_{\phi} + \frac{L_{\phi}}{\alpha} \right) \|\mathbf{x}^+ - \mathbf{x}\|.$$

Note that $\|\mathbf{x}^+ - \mathbf{x}\| = \alpha \|\mathcal{G}_{\alpha}(\mathbf{x})\|$, which completes the proof. \square

Proof of Proposition 9 By the definition of $\text{Prox}_R^{\phi}(\cdot)$, $x_i \in \partial R(\mathbf{x}_i^+) + \nabla \phi(\mathbf{x}_i^+)$, $i = 1, 2$. Since $\partial R(\cdot)$ is monotone, then $\langle \mathbf{x}_1 - \mathbf{x}_2 - (\nabla \phi(\mathbf{x}_1^+) - \nabla \phi(\mathbf{x}_2^+)), \mathbf{x}_1^+ - \mathbf{x}_2^+ \rangle \geq 0$. From the μ -strong convexity of ϕ , it follows that $\langle \mathbf{x}_1 - \mathbf{x}_2, \mathbf{x}_1^+ - \mathbf{x}_2^+ \rangle \geq \langle \nabla \phi(\mathbf{x}_1^+) - \nabla \phi(\mathbf{x}_2^+), \mathbf{x}_1^+ - \mathbf{x}_2^+ \rangle \geq \mu \|\mathbf{x}_1^+ - \mathbf{x}_2^+\|^2$. Therefore, $\|\mathbf{x}_1^+ - \mathbf{x}_2^+\| \leq \frac{1}{\mu} \|\mathbf{x}_1 - \mathbf{x}_2\|$. \square

Appendix B. Algorithmic stability analysis

Consider a dataset $\Xi = (\xi_1, \dots, \xi_n)$ and an algorithm \mathcal{A} . Let $\Phi_{\Xi}(\mathbf{x}) = F(\mathbf{x}, \Xi) + R(\mathbf{x})$, and let $\mathcal{A}(\Xi)$ be the output of the algorithm \mathcal{A} based on the dataset Ξ . Since the underlying distribution of ξ is unknown, the population risk can be decomposed into two parts:

$$\mathbb{E}_{\Xi, \mathcal{A}}[\Phi(\mathcal{A}(\Xi)) - \Phi(\mathbf{x}^*)] = \underbrace{\mathbb{E}_{\Xi, \mathcal{A}}[\Phi_{\Xi}(\mathcal{A}(\Xi)) - \Phi(\mathcal{A}(\Xi))]}_{\varepsilon_{\text{gen}}} + \underbrace{\mathbb{E}_{\Xi, \mathcal{A}}[\Phi_{\Xi}(\mathcal{A}(\Xi)) - \Phi(\mathbf{x}^*)]}_{\varepsilon_{\text{opt}}},$$

where \mathbf{x}^* is independent of Ξ and \mathcal{A} . The first term represents the expected generalization error, and the second term represents the optimization error. Our goal is to assess the expected generalization error ε_{gen} . For clarity, we focus on the simple case where $R(x) \equiv 0$, in which case $\Phi(\mathbf{x}) = F(\mathbf{x}) = \mathbb{E}_{\xi}[f(\mathbf{x}, \xi)]$. Consider any two data sets Ξ, Ξ' that differ by at most one example. As shown by (Bousquet and Elisseeff, 2002; Shalev-Shwartz et al., 2010; Hardt et al., 2016), the absolute expected generalization error $|\varepsilon_{\text{gen}}|$ can be bounded by ε if the algorithm is ε -uniform stable that is defined as follows:

Definition 31 *A randomized algorithm \mathcal{A} is ε -uniformly stable, if for any two datasets Ξ and Ξ' which differ by at most one example, it holds that*

$$\sup_z \mathbb{E}_{\mathcal{A}}[f(\mathcal{A}(\Xi), z) - f(\mathcal{A}(\Xi'), z)] \leq \varepsilon.$$

Lemma 32 *Let \mathcal{A} be ε -uniformly stable, then $|\varepsilon_{\text{gen}}| \leq \varepsilon$.*

In this section, we demonstrate that the polynomial kernel employed in our method (as discussed in Section 4) enhances algorithmic stability, particularly in high-dimensional scenarios such as training deep neural networks. Inspired by the ODE approach for Bregman gradient-type methods (Alvarez et al., 2004; Ding and Toh, 2024), the ODE corresponding to SBPG is given by:

$$\frac{d\nabla\phi(\mathbf{x}(t))}{dt} = -\nabla F(\mathbf{x}(t)),$$

which is equivalent to the following ODE:

$$\dot{\mathbf{x}}(t) = -[\nabla^2\phi(\mathbf{x}(t))]^{-1}\nabla F(\mathbf{x}(t)).$$

The discrete version of this ODE then leads to the Hessian preconditioned gradient method:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k[\nabla^2\phi(\mathbf{x}^k)]^{-1}\nabla F(\mathbf{x}^k).$$

This connection between the Bregman gradient method and the Hessian preconditioned gradient method inspires us to derive the following estimation on the expected generalization gap:

Theorem 33 *Given a dataset Ξ containing n samples. Let \mathcal{A} be the method (SBPG) where the last iterate is the output, and let \mathcal{E} be the event where the iterates generated by \mathcal{A} is contained within a compact set \mathcal{B} . Suppose f is differentiable and the kernel function ϕ is twice differentiable. Under the event \mathcal{E} , for any sufficiently small $\{\alpha_k\}$, we have*

$$|\mathbb{E}_{\Xi, \mathcal{A}}[\Phi_{\Xi}(\mathcal{A}(\Xi)) - \Phi(\mathcal{A}(\Xi))]| \leq \sum_{t=1}^k \exp\left(\left(1 - \frac{1}{n}\right) L_{\mathcal{B}} \sum_{i=t+1}^k \alpha_i\right) \frac{3\ell_{f, \mathcal{B}}^2 \ell_{\mathcal{B}} \alpha_t}{n} =: \varepsilon_{\text{gen}}^{\phi}, \quad (28)$$

where $L_{\mathcal{B}} = \sup_{\mathbf{x} \in \mathcal{B}, \mathbf{y} \in \mathcal{B}, \mathbf{z}} \frac{\|[\nabla^2 \phi(\mathbf{x})]^{-1} \nabla f(\mathbf{x}, \mathbf{z}) - [\nabla^2 \phi(\mathbf{y})]^{-1} \nabla f(\mathbf{y}, \mathbf{z})\|}{\|\mathbf{x} - \mathbf{y}\|}$, $\ell_{f, \mathcal{B}} = \sup_{\mathbf{x} \in \mathcal{B}, \mathbf{y} \in \mathcal{B}, \mathbf{z}} \frac{\|f(\mathbf{x}, \mathbf{z}) - f(\mathbf{y}, \mathbf{z})\|}{\|\mathbf{x} - \mathbf{y}\|}$, and $\bar{\ell}_{\mathcal{B}} = \sup_{\mathbf{x} \in \mathcal{B}} \|[\nabla^2 \phi(\mathbf{x})]^{-1}\|$.

Proof Given the dataset Ξ , recall the update scheme of SBPG:

$$\mathbf{x}_{k+1} = \nabla \phi^*(\nabla \phi(\mathbf{x}_k) - \alpha_k \nabla f(\mathbf{x}_k, \mathbf{z}_k)).$$

Noting that $\nabla \phi^*(\nabla \phi(\mathbf{x})) = \mathbf{x}$ and $\nabla^2 \phi^*(\nabla \phi(\mathbf{x})) = [\nabla^2 \phi(\mathbf{x})]^{-1}$, we have

$$\begin{aligned} \mathbf{x}_{k+1} &= \nabla \phi^*(\nabla \phi(\mathbf{x}_k) - \alpha_k \nabla f(\mathbf{x}_k, \mathbf{z}_k)) \\ &= \nabla \phi^*(\nabla \phi(\mathbf{x}_k)) - \alpha_k \nabla^2 \phi^*(\nabla \phi(\mathbf{x}_k)) \nabla f(\mathbf{x}_k, \mathbf{z}_k) + \mathcal{O}(\alpha_k^2) \\ &= \mathbf{x}_k - \alpha_k [\nabla^2 \phi(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k, \mathbf{z}_k) + \mathcal{O}(\alpha_k^2), \end{aligned}$$

where the second equality comes from Taylor's expansion. Let Ξ' be a dataset that differs from Ξ by one example. Define $\delta_{k+1} := \mathbb{E} [\|\mathbf{x}_{k+1} - \mathbf{x}'_{k+1}\|]$, where \mathbf{x}'_{k+1} corresponds to the dataset Ξ' . From the inequality above, we have

$$\begin{aligned} &\|\mathbf{x}_{k+1} - \mathbf{x}'_{k+1}\| \\ &= \|(\mathbf{x}_k - \alpha_k [\nabla^2 \phi(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k, \mathbf{z}_k)) - (\mathbf{x}'_k - \alpha_k [\nabla^2 \phi(\mathbf{x}'_k)]^{-1} \nabla f(\mathbf{x}'_k, \mathbf{z}'_k))\| + \mathcal{O}(\alpha_k^2) \\ &\leq \|\mathbf{x}_k - \mathbf{x}'_k\| + \alpha_k \|[\nabla^2 \phi(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k, \mathbf{z}_k) - [\nabla^2 \phi(\mathbf{x}'_k)]^{-1} \nabla f(\mathbf{x}'_k, \mathbf{z}'_k)\| + \mathcal{O}(\alpha_k^2). \end{aligned}$$

Since $\mathbf{x}_0 = \mathbf{x}'_0$, then we have

$$\begin{aligned} \delta_{k+1} &= \left(1 - \frac{1}{n}\right) (1 + \alpha_k L_{\mathcal{B}}) \delta_k + \frac{1}{n} (\delta_k + 2\ell_{f, \mathcal{B}} \bar{\ell}_{\mathcal{B}} \alpha_k) + C \alpha_k^2 \\ &\leq \left(1 + \alpha_k \left(1 - \frac{1}{n}\right) L_{\mathcal{B}}\right) \delta_k + \frac{3\ell_{f, \mathcal{B}} \bar{\ell}_{\mathcal{B}} \alpha_k}{n} \\ &\leq \exp\left(\alpha_k \left(1 - \frac{1}{n}\right) L_{\mathcal{B}}\right) \delta_k + \frac{3\ell_{f, \mathcal{B}} \bar{\ell}_{\mathcal{B}} \alpha_k}{n}. \end{aligned} \tag{29}$$

The first inequality comes from the fact that α_k is sufficiently small. Since $\delta_0 = 0$, by recursion (29), we have

$$\delta_k \leq \sum_{t=1}^k \exp\left(\left(1 - \frac{1}{n}\right) L_{\mathcal{B}} \sum_{i=t+1}^k \alpha_i\right) \frac{3\ell_{f, \mathcal{B}} \bar{\ell}_{\mathcal{B}} \alpha_t}{n}.$$

Finally, using the bound $\sup_{\mathbf{z}} \mathbb{E}_{\mathcal{A}} [f(\mathbf{x}_k, \mathbf{z}) - f(\mathbf{x}'_k, \mathbf{z})] \leq \ell_{f, \mathcal{B}} \delta_k$ and Lemma 32, we complete the proof. \blacksquare

Now, we make some remarks on Theorem 33.

Remark 34 1. When $\phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2$, we use the notation $\bar{L}_{\mathcal{B}} = \sup_{\mathbf{x} \in \mathcal{B}, \mathbf{y} \in \mathcal{B}, \mathbf{z}} \frac{\|\nabla f(\mathbf{x}, \mathbf{z}) - \nabla f(\mathbf{y}, \mathbf{z})\|}{\|\mathbf{x} - \mathbf{y}\|}$, $\bar{\ell}_{\mathcal{B}} = 1$. If $\alpha_k \leq \frac{c}{k}$ for some constant c , the right-hand side of (28) recovers the bound in (Hardt et al., 2016, Theorem 3.12).

2. If we choose the polynomial kernel from Proposition 27 with $\delta = 1$ and $r = 4$, i.e. $\phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2 + \frac{1}{4} \|\mathbf{x}\|^4$, then $\nabla^2 \phi(\mathbf{x})^{-1} = \frac{1}{1+\|\mathbf{x}\|^2} I - \frac{2\mathbf{x}\mathbf{x}^T}{(1+3\|\mathbf{x}\|^2)(1+\|\mathbf{x}\|^2)}$. Thus, we have $\ell_{\mathcal{B}} = \sup_{\mathbf{x} \in \mathcal{B}} \left\| \frac{1}{1+\|\mathbf{x}\|^2} I - \frac{2\mathbf{x}\mathbf{x}^T}{(1+3\|\mathbf{x}\|^2)(1+\|\mathbf{x}\|^2)} \right\| \leq 1$. When \mathcal{B} is a convex compact set, by the intermediate theorem, we have

$$L_{\mathcal{B}} = \ell_{f, \mathcal{B}} \sup_{\mathbf{x} \in \mathcal{B}} \left\{ \left\| D_{\nabla^2 \phi(\mathbf{x})^{-1}} \right\| \right\} + \ell_{\mathcal{B}} \bar{L}_{\mathcal{B}},$$

where $D_{\nabla^2 \phi(\mathbf{x})^{-1}}$ is the first order differential operator of $\nabla^2 \phi(\mathbf{x})^{-1}$. Furthermore, if $\text{dist}(0, \mathcal{B}) \geq M > 0$, then $\ell_{\mathcal{B}} \leq \frac{1}{1+M^2} < 1$. Moreover by some basic algebraic cal-

culations, we have $D_{\mathcal{B}} := \sup_{\mathbf{x} \in \mathcal{B}} \left\{ \left\| D_{\nabla^2 \phi(\mathbf{x})^{-1}} \right\| \right\} \leq \begin{cases} 6 \|\mathbf{x}\| (1 + 4 \|\mathbf{x}\|^2), & \|\mathbf{x}\| \leq 1, \\ \frac{12}{1+3\|\mathbf{x}\|^2}, & \|\mathbf{x}\| > 1. \end{cases}$

When M is sufficiently large, which usually occurs in high-dimensional scenarios, $D_{\mathcal{B}}$ becomes very small. Thus, $L_{\mathcal{B}} < \bar{L}_{\mathcal{B}}$. In summary, when a high-order polynomial kernel is employed, (SBPG) tends to have a better stability bound compared to standard SGD.

Appendix C. Additional Experimental Results

In this section, we provide further experimental results and assess the performance of MSBPG on more recent neural network architectures, including ConvNext (Liu et al., 2022) and the Vision Transformer (ViT) (Dosovitskiy, 2020). The networks are trained on the CIFAR-100 dataset for 200 epochs, utilizing a cosine annealing learning rate schedule (Loshchilov and Hutter, 2016), with a batch size of 128. The test accuracies for different optimizers, including MSBPG, SGD, Adam, and AdamW, are reported in Table 1. Notably, MSBPG achieved slightly higher test accuracy compared to the other methods. This advantage can be attributed to the Bregman proximity model used in our approach.

Method	MSBPG	SGD	Adam	AdamW
ViT Tiny	62.37	60.65	56.88	58.77
ViT Small	63.59	61.90	57.66	59.39
ConvNext Atto	75.98	74.46	73.76	75.29
ConvNext Femto	76.47	75.96	74.35	75.65

Table 1: Test accuracy (%) of different optimizers on CIFAR-100 dataset.

To further evaluate the performance of MSBPG, we applied it to Transformer-XL (Dai et al., 2019), a model designed for long-sequence tasks. We followed the official configuration to train the Transformer-XL-based model on the WikiText-103 dataset (Merity et al., 2016), a large-scale word-level language modeling benchmark that involves long-term dependencies. The performance of MSBPG, along with that of SGD, Adam, and AdamW, was measured by the test perplexity after 50,000 training steps. The results are presented in Table 2.

We also plot the gradient norm of the objective function, which indicates the stationarity of the minimization problem when the nonsmooth term is absent, as shown in Figure 11. From the figure, we observe that our method, MSBPG, successfully reaches a stationary point, similar to other optimization methods. It is important to highlight that our method

	MSBPG	SGD	Adam	AdamW
Transformer-XL	31.07	33.81	33.53	32.17

Table 2: Test perplexity (lower is better) for Transformer-XL-based model on the WikiText-103 dataset.

theoretically converges to the stationary point without requiring Lipschitz smoothness, while the convergence of traditional methods relies on this property. With careful tuning of the learning rate, methods such as SGD, Adam, and AdamW can also achieve stationary points, but they are more sensitive to learning rate choices compared to MSBPG. This sensitivity is further demonstrated experimentally in Figure 10, and is partially attributed to the absence of Lipschitz smoothness, which increases the sensitivity to the learning rate.

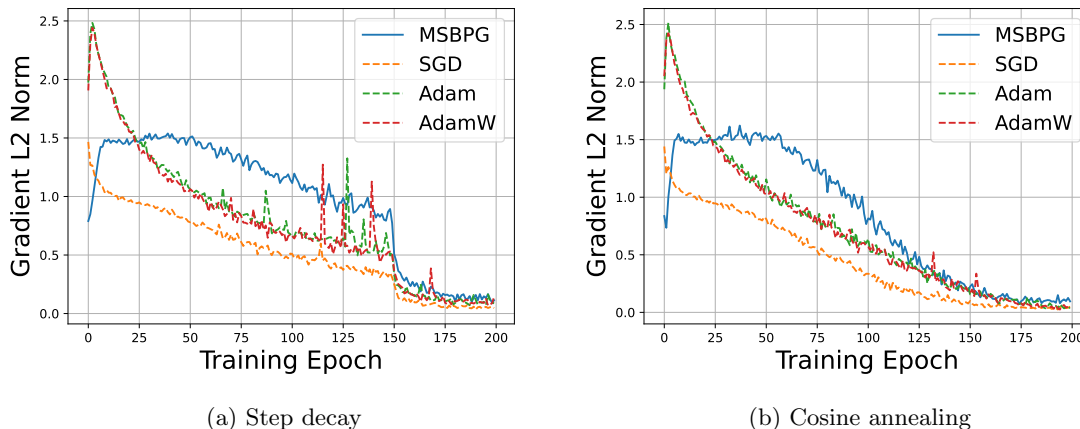


Figure 11: Gradient norm of the objective function. When the nonsmooth term is absent, the gradient norm can imply stationarity of the minimization problem.

Appendix D. Additional Figures

References

- Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. *Journal of Machine Learning Research*, 18(221):1–51, 2018.
- Felipe Alvarez, Jérôme Bolte, and Olivier Brahic. Hessian riemannian gradient flows in convex programming. *SIAM journal on control and optimization*, 43(2):477–501, 2004.
- Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth. Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 199(1):165–214, 2023.

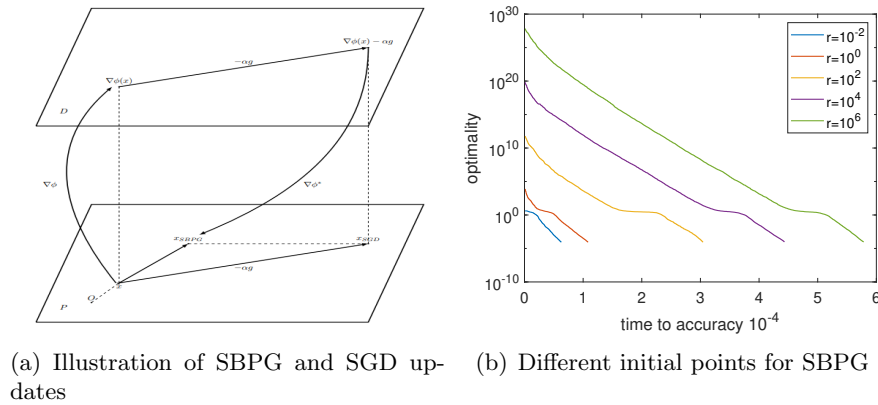


Figure 12: Figure (a) depicts SGD and SBPG updates. SBPG includes a "pull back" mechanism that prevents the point from moving excessively in any given direction. "P" and "D" refer to the primal and dual spaces, respectively, and these terms are commonly used in the mirror descent method literature (see, e.g., Bubeck et al. (2015); Nemirovskij and Yudin (1983)). Figure (b) illustrates the effect of choosing the initial point from a ball of radius r on the SBPG when r changes for the QIP example with $d = 100$ and $n = 5000$. All initial step sizes are set to 1×10^{-3} . As shown in Figure (b), even for an initial point that is far from the optimal point, SBPG can pull back the iterates to the optimal point.

Navid Azizan, Sahin Lale, and Babak Hassibi. Stochastic mirror descent on overparameterized nonlinear models: Convergence, implicit regularization, and generalization. *arXiv preprint arXiv:1906.03830*, 2019.

Heinz H Bauschke, Jérôme Bolte, and Marc Teboulle. A descent lemma beyond Lipschitz gradient continuity: first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.

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

Dimitri P Bertsekas. Incremental proximal methods for large scale convex optimization. *Mathematical Programming*, 129(2):163–195, 2011.

Dimitri P Bertsekas and John N Tsitsiklis. Gradient convergence in gradient methods with errors. *SIAM Journal on Optimization*, 10(3):627–642, 2000.

Pascal Bianchi. Ergodic convergence of a stochastic proximal point algorithm. *SIAM Journal on Optimization*, 26(4):2235–2260, 2016.

Jérôme Bolte, Shoham Sabach, Marc Teboulle, and Yakov Vaisbourd. First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems. *SIAM Journal on Optimization*, 28(3):2131–2151, 2018.

- Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002.
- Lev M Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–357, 2015.
- Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.
- Jinghui Chen, Dongruo Zhou, Yiqi Tang, Ziyang Yang, Yuan Cao, and Quanquan Gu. Closing the generalization gap of adaptive gradient methods in training deep neural networks. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3267–3275, 2021.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, 2019.
- Aaron Defazio and Léon Bottou. On the ineffectiveness of variance reduced optimization for deep learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Kuangyu Ding and Kim-Chuan Toh. Stochastic bregman subgradient methods for nonsmooth nonconvex optimization problems. *arXiv preprint arXiv:2404.17386*, 2024.
- Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Radu-Alexandru Dragomir, Alexandre d’Aspremont, and Jérôme Bolte. Quartic first-order methods for low-rank minimization. *Journal of Optimization Theory and Applications*, 189:341–363, 2021a.
- Radu Alexandru Dragomir, Mathieu Even, and Hadrien Hendrikx. Fast stochastic bregman gradient methods: Sharp analysis and variance reduction. In *International Conference on Machine Learning*, pages 2815–2825. PMLR, 2021b.
- John Duchi and Yoram Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10:2899–2934, 2009.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7), 2011.

- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2010.
- Rui Fu, Zuo Zhang, and Li Li. Using LSTM and GRU neural network methods for traffic flow prediction. In *2016 31st Youth Academic Annual Conference of Chinese Association of Automation (YAC)*, pages 324–328. IEEE, 2016.
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances in neural information processing systems*, 29, 2016.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-2):267–305, 2016.
- Bin Gu, Wenhan Xian, Zhouyuan Huo, Cheng Deng, and Heng Huang. A unified q-memorization framework for asynchronous stochastic optimization. *The Journal of Machine Learning Research*, 21(1):7761–7813, 2020.
- Filip Hanzely and Peter Richtárik. Fastest rates for stochastic mirror descent methods. *Computational Optimization and Applications*, 79:717–766, 2021.
- Filip Hanzely, Peter Richtarik, and Lin Xiao. Accelerated Bregman proximal gradient methods for relatively smooth convex optimization. *Computational Optimization and Applications*, 79:405–440, 2021.
- Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016.
- Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Convex analysis and minimization algorithms I: Fundamentals*, volume 305. Springer science & business media, 1993.
- Sepp Hochreiter. Untersuchungen zu dynamischen neuronalen netzen. *Diploma, Technische Universität München*, 91(1), 1991.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. pmlr, 2015.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Puya Latafat, Andreas Themelis, Masoud Ahookhosh, and Panagiotis Patrinos. Bregman finito/miso for nonconvex regularized finite sum minimization without lipschitz gradient continuity. *SIAM Journal on Optimization*, 32(3):2230–2262, 2022.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553): 436–444, 2015.
- Jason D Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I Jordan, and Benjamin Recht. First-order methods almost always avoid strict saddle points. *Mathematical programming*, 176:311–337, 2019.
- Yan Li, Caleb Ju, Ethan X Fang, and Tuo Zhao. Implicit regularization of bregman proximal point algorithm and mirror descent on separable data. *arXiv preprint arXiv:2108.06808*, 2021.
- Chih-Jen Lin, Ruby C Weng, and S Sathiya Keerthi. Trust region newton methods for large-scale logistic regression. In *Proceedings of the 24th international conference on Machine learning*, pages 561–568, 2007.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Haihao Lu. “Relative continuity” for non-Lipschitz nonsmooth convex optimization using stochastic (or deterministic) mirror descent. *INFORMS Journal on Optimization*, 1(4): 288–303, 2019.
- Haihao Lu, Robert M Freund, and Yurii Nesterov. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.

- D Russell Luke. Phase retrieval, what’s new. *SIAG/OPT Views and News*, 25(1):1–5, 2017.
- Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- Nikolay Manchev and Michael Spratling. Target propagation in recurrent neural networks. *The Journal of Machine Learning Research*, 21(1):250–282, 2020.
- Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. *Using Large Corpora*, 273, 1994.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- Arkadij Semenovič Nemirovskij and David Borisovich Yudin. Problem complexity and method efficiency in optimization. *Wiley Interscience*, 1983.
- Yu Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103:127–152, 2005.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- Andrew Y Ng. Feature selection, L1 vs. L2 regularization, and rotational invariance. In *Proceedings of the twenty-first International Conference on Machine Learning*, page 78, 2004.
- Ioannis Panageas, Georgios Piliouras, and Xiao Wang. First-order methods almost always avoid saddle points: The case of vanishing step-sizes. *Advances in Neural Information Processing Systems*, 32, 2019.
- Alasdair Paren, Leonard Berrada, Rudra PK Poudel, and M Pawan Kumar. A stochastic bundle method for interpolating networks. *Journal of Machine Learning Research*, 23: 1–57, 2022.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International Conference on Machine Learning*, pages 1310–1318. Pmlr, 2013.
- Andrei Patrascu and Ion Necoara. Nonasymptotic convergence of stochastic proximal point methods for constrained convex optimization. *The Journal of Machine Learning Research*, 18(1):7204–7245, 2017.
- Boris Teodorovich Polyak. Gradient methods for minimizing functionals. *Zhurnal vychislitel’noi matematiki i matematicheskoi fiziki*, 3(4):643–653, 1963.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.
- R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, 14(5):877–898, 1976.

- R Tyrrell Rockafellar. *Convex analysis*, volume 11. Princeton university press, 1997.
- R. Tyrrell Rockafellar and Roger J.-B. Wets. *Variational Analysis*. Springer Verlag, Heidelberg, Berlin, New York, 1998.
- Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *The Journal of Machine Learning Research*, 11:2635–2670, 2010.
- Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2021.
- Sima Siami-Namini, Neda Tavakoli, and Akbar Siami Namin. The performance of LSTM and BiLSTM in forecasting time series. In *2019 IEEE International Conference on Big Data*, pages 3285–3292. IEEE, 2019.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Haoyuan Sun, Kwangjun Ahn, Christos Thrampoulidis, and Navid Azizan. Mirror descent maximizes generalized margin and can be implemented efficiently. *Advances in Neural Information Processing Systems*, 35:31089–31101, 2022.
- Haoyuan Sun, Khashayar Gatmiry, Kwangjun Ahn, and Navid Azizan. A unified approach to controlling implicit regularization via mirror descent. *Journal of Machine Learning Research*, 24(393):1–58, 2023.
- Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, 18:1131–1198, 2018.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International Conference on Machine Learning*, pages 1139–1147. PMLR, 2013.
- Marc Teboulle. A simplified view of first order methods for optimization. *Mathematical Programming*, 170(1):67–96, 2018.
- Bokun Wang, Shiqian Ma, and Lingzhou Xue. Riemannian stochastic proximal gradient methods for nonsmooth optimization over the stiefel manifold. *The Journal of Machine Learning Research*, 23(1):4599–4631, 2022.
- Mengdi Wang, Ethan X Fang, and Han Liu. Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions. *Mathematical Programming*, 161:419–449, 2017.
- Qingsong Wang and Deren Han. A Bregman stochastic method for nonconvex nonsmooth problem beyond global Lipschitz gradient continuity. *Optimization Methods and Software*, 38(5):914–946, 2023.

- Lei Yang and Kim-Chuan Toh. Inexact Bregman Proximal Gradient Method and its Inertial Variant with Absolute and Relative Stopping Criteria. *arXiv preprint arXiv:2109.05690*, 2021.
- Lei Yang and Kim-Chuan Toh. Bregman proximal point algorithm revisited: A new inexact version and its inertial variant. *SIAM Journal on Optimization*, 32(3):1523–1554, 2022.
- Lei Yang, Jingjing Hu, and Kim-Chuan Toh. An inexact bregman proximal difference-of-convex algorithm with two types of relative stopping criteria. *arXiv preprint arXiv:2406.04646*, 2024.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Tong Zhang. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proceedings of the twenty-first International Conference on Machine Learning*, page 116, 2004.
- Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar C Tatikonda, Nicha Dvornek, Xenophon Papademetris, and James Duncan. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients. *Advances in Neural Information Processing Systems*, 33:18795–18806, 2020.