

Mixed Optimization for Smooth Functions

Authored by:

Rong Jin
Mehrdad Mahdavi
Lijun Zhang

Abstract

It is well known that the optimal convergence rate for stochastic optimization of smooth functions is $O(1/\sqrt{T})$, which is same as stochastic optimization of Lipschitz continuous convex functions. This is in contrast to optimizing smooth functions using full gradients, which yields a convergence rate of $O(1/T^2)$. In this work, we consider a new setup for optimizing smooth functions, termed as **Mixed Optimization**, which allows to access both a stochastic oracle and a full gradient oracle. Our goal is to significantly improve the convergence rate of stochastic optimization of smooth functions by having an additional small number of accesses to the full gradient oracle. We show that, with an $O(\ln T)$ calls to the full gradient oracle and an $O(T)$ calls to the stochastic oracle, the proposed mixed optimization algorithm is able to achieve an optimization error of $O(1/T)$.

1 Paper Body

Many machine learning algorithms follow the framework of empirical risk minimization, which often can be cast into the following generic optimization problem $\min_{w \in W} \sum_{i=1}^n g_i(w)$, where n is the number of training examples, $g_i(w)$ encodes the loss function related to the i th training example (x_i, y_i) , and W is a bounded convex domain that is introduced to regularize the solution $w \in W$ (i.e., the smaller the size of W , the stronger the regularization is). In this study, we focus on the learning problems for which the loss function $g_i(w)$ is smooth. Examples of smooth loss functions include least square with $g_i(w) = (y_i - w \cdot x_i)^2$ and logistic regression with $g_i(w) = \log(1 + \exp(-y_i \cdot w \cdot x_i))$. Since the regularization is enforced through the restricted domain W , we did not introduce a $\|w\|_2^2/2$ regularizer into the optimization problem and as a result, we do not assume the loss function to be strongly convex. We note that a small $\|w\|_2^2$ regularizer does NOT improve the convergence rate of stochastic optimization. More specifically, the convergence rate for stochastically optimizing a $\|w\|_2^2$ regularized loss function remains as $O(1/\sqrt{T})$ when ?

$= O(1/\sqrt{T})$ [11, Theorem 1], a scenario that is often encountered in real-world applications. A preliminary approach for solving the optimization problem in (1) is the batch gradient descent (GD) algorithm [16]. It starts with some initial point, and iteratively updates the solution using the equation $\mathbf{w}_{t+1} = \Pi_W(\mathbf{w}_t - \eta \nabla G(\mathbf{w}_t))$, where $\Pi_W(\cdot)$ is the orthogonal projection onto the convex domain W . It has been shown that for smooth objective functions, the convergence rate of standard GD is $O(1/\sqrt{T})$ [16], and can be improved to $O(1/T^2)$ by an accelerated GD algorithm [15, 16, 18]. The main shortcoming of GD method is its high cost in computing the full gradient $\nabla G(\mathbf{w}_t)$ when the number of training examples is large. Stochastic gradient descent (SGD) [3, 13, 21] alleviates this limitation of GD by sampling one (or a small set of) examples and computing a stochastic (sub)gradient at each iteration based on the sampled examples. Since the computational cost of SGD per iteration is independent of the size of the data (i.e., n), it is usually appealing for largescale learning and optimization. While SGD enjoys a high computational efficiency per iteration, it suffers from a slow convergence rate for optimizing smooth functions. It has been shown in [14] that the effect of the stochastic noise

Full (GD)	$O(\sqrt{T})$	$O(1)$	$O(1)$	$O(1)$
Setting	Lipschitz			
Convergence				
Smooth				
	$1/\sqrt{T}$	$1/T^2$		
	0			
Stochastic (SGD)	$O(\sqrt{T})$	$O(1)$	$O(1)$	$O(1)$
	$O(\sqrt{T})$			
Convergence				
	$1/\sqrt{T}$			
	$1/\sqrt{T}$			
	0			
Mixed Optimization	$O(\sqrt{T})$	$O(1)$	$O(1)$	$O(1)$
Convergence	$1/\sqrt{T}$			
	$1/\sqrt{T}$			
	$\log T$			

Table 1: The convergence rate (O), number of calls to stochastic oracle (O_s), and number of calls to full gradient oracle (O_f) for optimizing Lipschitz continuous and smooth convex functions, using full GD, SGD, and mixed optimization methods, measured in the number of iterations T . $\sqrt{\cdot}$ cannot be decreased with a better rate than $O(1/\sqrt{T})$ which is significantly worse than GD that uses the full gradients for updating the solutions and this limitation is also valid when the target function is smooth. In addition, as we can see from Table 1, for general Lipschitz functions, SGD exhibits the same convergence rate as that for the smooth functions, implying that smoothness is essentially not very useful and can not be exploited in stochastic optimization. The slow convergence rate for stochastically optimizing smooth loss functions is mostly due to the variance in stochastic gradients: unlike the full gradient case where the norm of a gradient

approaches to zero when the solution is approaching to the optimal solution, in stochastic optimization, the norm of a stochastic gradient is constant even when the solution is close to the optimal solution. It is the variance in stochastic gradients that makes the convergence rate $O(1/\sqrt{T})$ unimprovable in smooth setting [14, 1]. In this study, we are interested in designing an efficient algorithm that is in the same spirit of SGD but can effectively leverage the smoothness of the loss function to achieve a significantly faster convergence rate. To this end, we consider a new setup for optimization that allows us to interplay between stochastic and deterministic gradient descent methods. In particular, we assume that the optimization algorithm has an access to two oracles: \mathcal{O}_s A stochastic oracle that returns the loss function $g_i(w)$ and its gradient based on the sampled training example (x_i, y_i) , and \mathcal{O}_f A full gradient oracle that returns the gradient $\nabla G(w)$ for any given solution $w \in W$. We refer to this new setting as mixed optimization in order to distinguish it from both stochastic and full gradient optimization models. The key question we examined in this study is: Is it possible to improve the convergence rate for stochastic optimization of smooth functions by having a small number of calls to the full gradient oracle \mathcal{O}_f ? We give an affirmative answer to this question. We show that with an additional $O(\ln T)$ accesses to the full gradient oracle \mathcal{O}_f , the proposed algorithm, referred to as MIXED GRAD, can improve the convergence rate for stochastic optimization of smooth functions to $O(1/T)$, the same rate for stochastically optimizing a strongly convex function [11, 19, 23]. MIXED GRAD builds off on multistage methods [11] and operates in epochs, but involves novel ingredients so as to obtain an $O(1/T)$ rate for smooth losses. In particular, we form a sequence of strongly convex objective functions to be optimized at each epoch and decrease the amount of regularization and shrink the domain as the algorithm proceeds. The full gradient oracle \mathcal{O}_f is only called at the beginning of each epoch. Finally, we would like to distinguish mixed optimization from hybrid methods that use growing sample-sizes as optimization method proceeds to gradually transform the iterates into the full gradient method [9] and batch gradient with varying sample sizes [6], which unfortunately make the iterations to be dependent to the sample size n as opposed to SGD. In contrast, MIXED GRAD is as an alternation of deterministic and stochastic gradient steps, with different of frequencies for each type of steps. Our result for mixed optimization is useful for the scenario when the full gradient of the objective function can be computed relatively efficient although it is still significantly more expensive than computing a stochastic gradient. An example of such a scenario is distributed computing where the computation of full gradients can be speeded up by having it run in parallel on many machines with each machine containing a relatively small subset of the entire training data. Of course, the latency due to the communication between machines will result in an additional cost for computing the full gradient in a distributed fashion. Outline The rest of this paper is organized as follows. We begin in Section 2 by briefly reviewing the literature on deterministic and stochastic optimization. In Section 3, we introduce the necessary definitions and discuss the assumptions that underlie our analysis. Section 4 describes the MIXED GRAD algorithm and states the main result

on its convergence rate. The proof of main result is given in Section 5. Finally, Section 6 concludes the paper and discusses few open questions. ¹

The convergence rate can be improved to $O(1/T)$ when the structure of the objective function is provided. We note that the stochastic oracle assumed in our study is slightly stronger than the stochastic gradient oracle as it returns the sampled function instead of the stochastic gradient. ²

²

²

More Related Work

Deterministic Smooth Optimization The convergence rate of gradient based methods usually depends on the analytical properties of the objective function to be optimized. When the objective function is strongly convex and smooth, it is well known that a simple GD method can achieve a linear convergence rate [5]. For a μ -non-smooth Lipschitz-continuous function, the optimal rate for μ the first order method is only $O(1/\sqrt{T})$ [16]. Although $O(1/\sqrt{T})$ rate is not improvable in general, several recent studies are able to improve this rate to $O(1/T)$ by exploiting the special structure of the objective function [18, 17]. In the full gradient based convex optimization, smoothness is a highly desirable property. It has been shown that a simple GD achieves a convergence rate of $O(1/T)$ when the objective function is smooth, which is further can be improved to $O(1/T^2)$ by using the accelerated gradient methods [15, 18, 16]. **Stochastic Smooth Optimization** Unlike the optimization methods based on full gradients, the smoothness assumption was μ not exploited by most stochastic optimization methods. In fact, it was shown in [14] that the $O(1/\sqrt{T})$ convergence rate for stochastic optimization cannot be improved even when the objective function is smooth. This classical result is further confirmed by the recent studies of composite bounds for the first order optimization methods [2, 12]. The smoothness of the objective function is exploited extensively in mini-batch stochastic optimization [7, 8], where the goal is not to improve the convergence rate but to reduce the variance in stochastic gradients and consequently the number of times for updating the solutions [24]. We finally note that the smoothness assumption coupled with the strong convexity of function is beneficial in stochastic setting and yields a geometric convergence in expectation using Stochastic Average Gradient (SAG) and Stochastic Dual Coordinate Ascent (SDCA) algorithms proposed in [20] and [22], respectively.

3 Preliminaries We use bold-face letters to denote vectors. For any two vectors $w, w' \in W$, we denote by $\langle w, w' \rangle$ the inner product between w and w' . Throughout this paper, we only consider the ℓ_2 -norm. We assume the objective function $G(w)$ defined in (1) to be the average of n convex loss functions. The same assumption was made in [20, 22]. We assume that $G(w)$ is minimized at some $w^* \in W$. Without loss of generality, we assume that $W \subset \mathcal{B}(R)$, a ball of radius R . Besides convexity of individual functions, we will also assume that each $g_i(w)$ is μ -smooth as formally defined below [16]. **Definition 1 (Smoothness).** A differentiable loss function $f(w)$ is said to be μ -smooth with respect to a norm $\|\cdot\|$, if it holds that $\|f(w) - f(w')\| \leq \mu \|w - w'\|$, $w, w' \in W$. ² The smoothness assumption also implies that

$\|f(w) - f(w')\| \leq L\|w - w'\|$ which is equivalent to $f(w)$ being L -Lipschitz continuous. In stochastic first-order optimization setting, instead of having direct access to $G(w)$, we only have access to a stochastic gradient oracle, which given a solution $w \in W$, returns the gradient $g_i(w)$ where i is sampled uniformly at random from $\{1, 2, \dots, n\}$. The goal of stochastic optimization is to find $w \in W$ such that the optimization error, $\|G(w) - G(w^*)\|$, is as small as possible. In the mixed optimization model considered in this study, we first relax the stochastic oracle O_s by assuming that it will return a randomly sampled loss function $g_i(w)$, instead of the gradient $g_i(w)$ for a given solution w . Second, we assume that the learner also has an access to the full gradient oracle O_f . Our goal is to significantly improve the convergence rate of stochastic gradient descent (SGD) by making a small number of calls to the full gradient oracle O_f . In particular, we show that by having only $O(\log T)$ accesses to the full gradient oracle and $O(T)$ accesses to the stochastic oracle, we can tolerate the noise in stochastic gradients and attain an $O(1/T)$ convergence rate for optimizing smooth functions. The audience may feel that this relaxation of stochastic oracle could provide significantly more information, and second order methods such as Online Newton [10] may be applied to achieve $O(1/T)$ convergence. We note (i) the proposed algorithm is a first order method, and (ii) although the Online Newton method yields a regret bound of $O(1/T)$, its convergence rate for optimization can be as low as $O(1/\sqrt{T})$ due to the concentration bound for Martingales. In addition, the Online Newton method is only applicable to exponential concave function, not any smooth loss function.

3

Algorithm 1 MIXED GRAD Input: step size η , domain size γ , the number of iterations T for the first epoch, the number of epoches m , regularization parameter λ , and shrinking parameter β . 1: Initialize w . 2: for $k = 1, \dots, m$ do 3: Construct the domain $W_k = \{w : w + \eta \sum_{i=1}^k g_i(w) \in W, \|w\| \leq \gamma\}$. 4: Call the full gradient oracle O_f for $G(w_k) = \frac{1}{n} \sum_{i=1}^n g_i(w_k)$. 5: Compute $g_k = G(w_k)$. 6: Initialize $w_{k+1} = w_k$. 7: for $t = 1, \dots, T_k$ do 8: Call stochastic oracle O_s to return a randomly selected loss function $g_{i_t}(w)$. 9: Compute the stochastic gradient as $g_{i_t}(w_k)$. 10: Update the solution by $w_{k+1} = w_k + \eta g_{i_t}(w_k)$. 11: end for 12: Set $w_{k+1} = w_k + \eta \sum_{t=1}^{T_k} g_{i_t}(w_k)$. 13: Set $\gamma_{k+1} = \beta \gamma_k$, $\eta_{k+1} = \eta_k / \beta$, and $T_{k+1} = \lfloor \beta T_k \rfloor$. 14: end for 15: Return w .

The analysis of the proposed algorithm relies on the strong convexity of intermediate loss functions introduced to facilitate the optimization as given below. **Definition 2 (Strong convexity).** A function $f(w)$ is said to be μ -strongly convex w.r.t a norm $\|\cdot\|$, if there exists a constant $\mu > 0$ (often called the modulus of strong convexity) such that it holds $f(w) \geq f(w') + \mu\|w - w'\|^2$, $\forall w, w' \in W$.

4

Mixed Stochastic/Deterministic Gradient Descent

We now turn to describe the proposed mixed optimization algorithm and state its convergence rate. The detailed steps of MIXED GRAD algorithm are shown in Algorithm 1. It follows the epoch gradient descent algorithm proposed in [11] for stochastically minimizing strongly convex functions and divides the optimization process into m epochs, but involves novel ingredients so as to obtain an $O(1/T)$ convergence rate. The key idea is to introduce a λ_k regularizer into the objective function to make it strongly convex, and gradually reduce the amount of regularization over the epochs. We also shrink the domain as the algorithm proceeds. We note that reducing the amount of regularization over time is closely-related to the classic proximal-point algorithms. Throughout the paper, we will use the subscript for the index of each epoch, and the superscript for the index of iterations within each epoch. Below, we describe the key idea behind MIXED GRAD. w^k be the solution obtained before the k th epoch, which is initialized to be 0 for the first epoch. Let w^k , resulting in the following. Instead of searching for w^* at the k th epoch, our goal is to find w^* by solving the optimization problem for the k th epoch $\min_{w \in W^k} F_k(w) = \min_{w \in W^k} \{ \sum_{i=1}^n g_i(w) + \frac{\lambda_k}{2} \|w\|^2 \}$ (2)

where W^k specifies the domain size of w and λ_k is the regularization parameter introduced at the k th epoch. By introducing the λ_k regularizer, the objective function in (2) becomes strongly convex, making it possible to exploit the technique for stochastic optimization of strongly convex function in order to improve the convergence rate. The domain size W^k and the regularization parameter λ_k are initialized to be $W^1 \subseteq W$ and $\lambda_1 > 0$, respectively, and are reduced by a constant factor $\beta \in (0, 1)$ every epoch, i.e., $\lambda_k = \lambda_1 / \beta^{k-1}$ and $W^k = \beta^{k-1} W^1$. By removing the constant term $\frac{\lambda_k}{2} \|w\|^2$ from the objective function in (2), we obtain the following optimization problem for the k th epoch $\min_{w \in W^k} F_k(w) = \min_{w \in W^k} \{ \sum_{i=1}^n g_i(w) + \frac{\lambda_k}{2} \|w\|^2 \}$ (3)

where $W^k = \{w : w + w^k \in W, w^k \in W^k\}$. We rewrite the objective function $F_k(w)$ as $\min_{w \in W^k} F_k(w) = \min_{w \in W^k} \{ \sum_{i=1}^n g_i(w) + \frac{\lambda_k}{2} \|w\|^2 \}$. By removing the constant term $\frac{\lambda_k}{2} \|w\|^2$ from the objective function in (2), we obtain the following optimization problem for the k th epoch $\min_{w \in W^k} F_k(w) = \min_{w \in W^k} \{ \sum_{i=1}^n g_i(w) + \frac{\lambda_k}{2} \|w\|^2 \}$ (3)

= where

(4)

$\lambda_k = \lambda_1 / \beta^{k-1}$ and $g_{bik}(w) = g_i(w + w^k) - \frac{\lambda_k}{2} \|w\|^2$. $g_i(w) = g_i(w + w^k) - \frac{\lambda_k}{2} \|w\|^2$

The main reason for using $g_{bik}(w)$ instead of $g_i(w)$ is to tolerate the variance in the stochastic gradients. To see this, from the smoothness assumption of $g_i(w)$ we obtain the following inequality for the norm of $g_{bik}(w)$ as:

$\|g_{bik}(w)\| \leq \|g_i(w + w^k) - g_i(w)\|$

As a result, since λ_k and W^k shrinks over epochs, then λ_k will approach to zero over epochs and consequentially $g_{bik}(w)$ approaches to zero, which allows us to effectively control the variance in stochastic gradients, a key to improving the convergence of stochastic optimization for smooth functions to $O(1/T)$. Using $F_k(w)$ in

[illegible]

where in the last step holds under the condition $\epsilon \leq 2$. By combining above inequalities, we obtain $\sum_{i=1}^n \mathbb{E} \langle \nabla F_m(w^i), w^i - w^* \rangle \leq \sum_{i=1}^n \mathbb{E} \langle \nabla F_m(w^i), w^i - w^* \rangle + 3m+1 \mathbb{E} \langle \nabla F_m(w^i), w^i - w^* \rangle$. $\nabla F_m(w)$ minimizes $F_m(w)$, for any w . Our final goal is to relate $F_m(w)$ to $\min_w G(w)$. Since $w^* \in \arg \min G(w)$, we have $\langle \nabla F_m(w^*), w^* - w^* \rangle = 0$. $F_m(w^*) \leq F_m(w^*) = (6) \mathbb{E} \langle \nabla F_m(w^*), w^* - w^* \rangle + m \mathbb{E} \langle \nabla F_m(w^*), w^* - w^* \rangle$. To this end, after the first m epochs, the key to bound $-F(w^*) \leq G(w^*)$ is to bound $\mathbb{E} \langle \nabla F_m(w^*), w^* - w^* \rangle$. w^1, w^2, \dots be the sequence of solutions epochs, we run Algorithm 1 with full gradients. Let w^k generated by Algorithm 1 after the first m epochs. For this sequence of solutions, Theorem 2 e k ϵ^k for any ϵ will hold deterministically as we deploy the full gradient for updating, i.e., $\nabla F_m(w^k)$. Since we reduce ϵ^k exponentially, ϵ^k will approach to zero and therefore the sequence $\{w^k\}$ will converge to w^* , one of the optimal solutions that minimize $G(w)$. Since w^* is the k ϵ^k for any $k \geq m+1$, we have limit of sequence $\{w^k\}$ and $\nabla F_m(w^*) \leq \mathbb{E} \langle \nabla F_m(w^*), w^* - w^* \rangle = 0$.

6

where the last step follows from the condition $\sum_{i=1}^n \alpha_i = 1$. Thus, $\sum_{i=1}^n \alpha_i \log \frac{1}{\alpha_i} = \sum_{i=1}^n \alpha_i \log \frac{1}{\alpha_i} + \sum_{i=1}^n \alpha_i \log \frac{1}{\alpha_i} = 2 \sum_{i=1}^n \alpha_i \log \frac{1}{\alpha_i} = 2H(\alpha)$.

$$\sum_{n=1}^N \sum_{i=1}^n \left(\frac{1}{2} \sum_{j=1}^{n-i+1} \left(\frac{1}{2} \sum_{k=1}^{n-i-j+1} g_i(w_k) + 2m_{11}^2 + ? \right) g_m(w_j) + 2m_{11}^2 \right)$$

(7)

By combining the bounds in (6) and (7), we have, with a probability $1 - 2m^{-2}$, $\sum_{i=1}^n |g_i(w) - g_i(w^*)| \leq O(1/T) \sum_{i=1}^n \sum_{j=1}^m |w_{ij} - w_{ij}^*|$ where $T = T_1$.

$$m \geq 1 \quad ? \quad k=0$$

?

$2k$

$$(\quad) T1 \text{ ? } 2m \text{ ? } 1 \text{ T1 } 2m \text{ ? } ? . = 2 \text{ ? } ? 1 \text{ 3}$$

We complete the proof by plugging in the stated values for β , β_1 and β_2 .

5.1

Proof of Theorem 2

For the convenience of discussion, we drop the subscript k for epoch just to simplify our notation. $\theta^k = w^k$ be the solution obtained before the start of Let $\theta^k = \theta^k$, $T = T^k$, $\theta = \theta^k$, $g = g^k$. Let $w^{k+1} = w^k$ be the solution obtained after running through the k th epoch. We the epoch k , and let

$$= (9)$$

$$\begin{aligned}
 & \quad b^? \cdot 2b^? \cdot 2^? w^? \cdot w^? w^{T+1} \cdot w^? \cdot 2b^? \cdot F(w^?) \cdot F(w^? + b^? \cdot \text{git}(w^?) + w^? \cdot g, w^? w^{T+1} \cdot 2^? \cdot 2^? \cdot 2^? \cdot b^? w^? b^? w^? b^? t), w^? w^? b^? \cdot 2^? \cdot b^? b^? \cdot w^? \cdot w^? b^? + b^? b^? \cdot F(b^?) + F(w^? b^? + F(\text{git}(w^? \cdot \text{git}(w^?) + b^? \cdot \text{git}(w^? \cdot 1 = 0, \text{ we have } \\
 & \text{By adding the inequality in Lemma 1 over all iterations, using the fact } w^? T \cdot b^? \cdot 2b^? \cdot 2^? w^? w^{T+1} \cdot w^? b^? \cdot F(w^?) \cdot F(w^? \cdot g, w^{T+1} \cdot 2^? \cdot 2^? \cdot t=1 +
 \end{aligned}$$

```

+
T ??
:=BT

```

using the fact $F(0) \leq F(w) + 2w$ and $\max(w, wT + 1) \leq w$, we have $b \leq F(wT + 1) \leq F(0) + 2(wT + 1) \leq F(wT + 1) + 2w$ and therefore $(T + 1) \leq 2b \leq F(wT) + F(w) + AT + BT + CT$. \square

The following lemmas bound AT , BT and CT . Lemma 2. For AT defined above we have $\text{AT} \leq 6T^2$. The following lemma upper bounds BT and CT . The proof is based on the Bernstein's inequality for Martingales [4] and is given in the Appendix. Lemma 3. With a probability $1 - \epsilon$, we have $\text{BT} \leq 2\ln(1/\epsilon) + 2T \ln(1/\epsilon)$ and $\text{CT} \leq 2\ln(1/\epsilon) + 2T \ln(1/\epsilon)$. Using Lemmas 2 and 3, by substituting the uppers bounds for AT , BT , and CT in (10), with a probability $1 - \epsilon$, we obtain $(T+1)^{-1} \sum_{t=1}^T F(w_t) \leq F(w) + 6T^2 + 3\ln(1/\epsilon) + 3T \ln(1/\epsilon)$. By choosing $\epsilon = 1/(2T^3)$, we have $(T+1)^{-1} \sum_{t=1}^T F(w_t) \leq F(w) + e^{-1/T^3}$. By using the fact $w_t \geq 3 \ln[1/\epsilon] - 2.5T^2$, we have and using the fact $w_t \geq 3 \ln[1/\epsilon] - 2.5T^2$, we have and using the fact $w_t \geq 3 \ln[1/\epsilon] - 2.5T^2$.

$\mathbb{E} F(w)$, and $\mathbb{E} = \mathbb{E}_w$. $T+1 \leq T+1$ Thus, when $T \geq \lceil 300 \frac{8}{\epsilon} \frac{2 \ln 1/\delta}{\epsilon} \rceil$, we have, with a probability $1 - 2\delta$, $2 \leq b \leq 2$, and $-\mathbb{E}(w) \leq \mathbb{E}(wb) - \frac{1}{4} \frac{2}{\epsilon} \leq (11) \frac{4}{\epsilon} \frac{2}{\epsilon} \leq \mathbb{E}_w b \leq \frac{1}{\epsilon} b \leq \frac{1}{\epsilon}$ to \mathbb{E}_w The next lemma relates $\mathbb{E}_w b \leq \frac{1}{\epsilon} \leq \mathbb{E}_w b \leq \frac{1}{\epsilon}$. Lemma 4. We have $\mathbb{E}_w b \leq \frac{1}{\epsilon} \leq \mathbb{E}_w b \leq \frac{1}{\epsilon}$. Combining the bound in (11) with Lemma 4, we have \mathbb{E}_w

6

Conclusions and Open Questions

We presented a new paradigm of optimization, termed as mixed optimization, that aims to improve the convergence rate of stochastic optimization by making a small number of calls to the full gradient oracle. We proposed the MIXED GRAD algorithm and showed that it is able to achieve an $O(1/T)$ convergence rate by accessing stochastic and full gradient oracles for $O(T)$ and $O(\log T)$ times, respectively. We showed that the MIXED GRAD algorithm is able to exploit the smoothness of the function, which is believed to be not very useful in stochastic optimization. In the future, we would like to examine the optimality of our algorithm, namely if it is possible to achieve a better convergence rate for stochastic optimization of smooth functions using $O(\ln T)$ accesses to the full gradient oracle. Furthermore, to alleviate the computational cost caused by $O(\log T)$ accesses to the full gradient oracle, it would be interesting to empirically evaluate the proposed algorithm in a distributed framework by distributing the individual functions among processors to parallelize the full gradient computation at the beginning of each epoch which requires $O(\log T)$ communications between the processors in total. Lastly, it is very interesting to check whether an $O(1/T^2)$ rate could be achieved by an accelerated method in the mixed optimization scenario, and whether linear convergence rates could be achieved in the strongly-convex case. Acknowledgments. The authors would like to thank the anonymous reviewers for their helpful and insightful comments. This work was supported in part by ONR Award N000141210431 and NSF (IIS-1251031).

8

2 References

- [1] A. Agarwal, P. L. Bartlett, P. D. Ravikumar, and M. J. Wainwright. Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization. *IEEE Transactions on Information Theory*, 58(5):3235–3249, 2012.
- [2] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Oper. Res. Lett.*, 31(3):167–175, 2003.
- [3] L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *NIPS*, pages 161–168, 2008.
- [4] S. Boucheron, G. Lugosi, and O. Bousquet. Concentration inequalities. In *Advanced Lectures on Machine Learning*, pages 208–240, 2003.
- [5] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [6] R. H. Byrd, G. M. Chin, J. Nocedal, and Y. Wu. Sample size selection in optimization methods for machine learning. *Mathematical programming*, 134(1):127–155, 2012.
- [7] A. Cotter, O. Shamir, N. Srebro,

and K. Sridharan. Better mini-batch algorithms via accelerated gradient methods. In NIPS, pages 1647?1655, 2011. [8] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao. Optimal distributed online prediction using mini-batches. The Journal of Machine Learning Research, 13:165?202, 2012. [9] M. P. Friedlander and M. Schmidt. Hybrid deterministic-stochastic methods for data fitting. SIAM Journal on Scientific Computing, 34(3):A1380?A1405, 2012. [10] E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization. Machine Learning, 69(2-3):169?192, 2007. [11] E. Hazan and S. Kale. Beyond the regret minimization barrier: an optimal algorithm for stochastic strongly-convex optimization. Journal of Machine Learning Research - Proceedings Track, 19:421?436, 2011. [12] Q. Lin, X. Chen, and J. Pena. A smoothing stochastic gradient method for composite optimization. arXiv preprint arXiv:1008.5204, 2010. [13] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. SIAM J. on Optimization, 19:1574?1609, 2009. [14] A. S. Nemirovsky and D. B. Yudin. Problem complexity and method efficiency in optimization. 1983. [15] Y. Nesterov. A method of solving a convex programming problem with convergence rate $o(1/k^2)$. In Soviet Mathematics Doklady, volume 27, pages 372?376, 1983. [16] Y. Nesterov. Introductory Lectures on Convex Optimization: A Basic Course. Kluwer Academic Publishers, 2004. [17] Y. Nesterov. Excessive gap technique in nonsmooth convex minimization. SIAM Journal on Optimization, 16(1):235?249, 2005. [18] Y. Nesterov. Smooth minimization of non-smooth functions. Math. Program., 103(1):127? 152, 2005. [19] A. Rakhlin, O. Shamir, and K. Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. In ICML, 2012. [20] N. L. Roux, M. W. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In NIPS, pages 2672?2680, 2012. [21] S. Shalev-Shwartz, Y. Singer, and N. Srebro. Pegasos: Primal estimated sub-gradient solver for svm. In ICML, pages 807?814, 2007. [22] S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. JMLR, 14:567599, 2013. [23] O. Shamir and T. Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. ICML, 2013. [24] L. Zhang, T. Yang, R. Jin, and X. He. $O(\log t)$ projections for stochastic optimization of smooth and strongly convex functions. ICML, 2013.