
Regret Bounds without Lipschitz Continuity: Online Learning with Relative-Lipschitz Losses

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In online convex optimization (OCO), Lipschitz continuity of the functions is
2 commonly assumed in order to obtain sublinear regret. Moreover, many algorithms
3 have only logarithmic regret when these functions are also strongly convex. Re-
4 cently, researchers from convex optimization proposed the notions of “relative
5 Lipschitz continuity” and “relative strong convexity”. Both of the notions are
6 generalizations of their classical counterparts. It has been shown that subgradient
7 methods in the relative setting have performance analogous to their performance in
8 the classical setting.

9 In this work, we consider OCO for relative Lipschitz and relative strongly convex
10 functions. We extend the known regret bounds for classical OCO algorithms to the
11 relative setting. Specifically, we show regret bounds for the follow the regularized
12 leader algorithms and a variant of online mirror descent. Due to the generality
13 of these methods, these results yield regret bounds for a wide variety of OCO
14 algorithms. Furthermore, we further extend the results to algorithms with extra
15 regularization such as regularized dual averaging.

16 1 Introduction

17 In online convex optimization (OCO), at each of many rounds a player has to pick a point from a
18 convex set while an adversary chooses a convex function that penalizes the player’s choice. More
19 precisely, in each round $t \in \mathbb{N}$, the player picks a point x_t from a fixed convex set $\mathcal{X} \subseteq \mathbb{R}^n$ and an
20 adversary picks a convex function f_t depending on x_t . At the end of the round, the player suffers a
21 loss of $f_t(x_t)$. Besides modeling a wide range of online learning problems [Shalev-Shwartz, 2011],
22 algorithms for OCO are often used in batch optimization problems due to their low computational
23 cost per iteration. For example, the widely used stochastic gradient descent (SGD) algorithm can be
24 viewed as a special case of online gradient descent [Hazan, 2016, Chapter 3] and AdaGrad [Duchi
25 et al., 2011] is a foundational adaptive gradient descent method originally proposed in the OCO
26 setting. The performance measure usually used for OCO algorithms is the *regret*. It is the difference
27 between the cost incurred to the player and a comparison point $z \in \mathcal{X} \subseteq \mathbb{R}^n$ (usually with minimum
28 cumulative loss), that is to say,

$$\text{Regret}_T(z) := \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(z).$$

29 Classical results show that if the cost functions are Lipschitz continuous, then there are algorithms
30 which suffer at most $O(\sqrt{T})$ regret in T rounds [Zinkevich, 2003]. Additionally, if the cost functions
31 are strongly convex, there are algorithms that suffer at most $O(\log T)$ regret in T rounds [Hazan
32 et al., 2007]). However, not all loss functions that appear in applications, such as in inverse Poisson

33 problems [Antonakopoulos et al., 2020] and support vector machines training [Lu, 2019], satisfy
34 these conditions on the entire feasible set.

35 Recently, there has been a line of work investigating the performance of optimization methods beyond
36 conventional assumptions [Bauschke et al., 2017, Lu et al., 2018, Lu, 2019]. Intriguingly, much of
37 this line of work proposes relaxed assumptions under which classical algorithms enjoy convergence
38 rates similar to the ones from the classical setting.

39 In particular, Lu [2019] proposed the notion of relative Lipschitz-continuity and showed how mirror
40 descent (with properly chosen regularizer/mirror map) converges at a rate of $O(1/\sqrt{T})$ in T iter-
41 ations for non-smooth relative Lipschitz-continuous functions. Furthermore, they show a $O(1/T)$
42 convergence rate when the function is also relatively strongly-convex (a notion proposed by Lu et al.
43 [2018]). Although the former result can be translated to a $O(\sqrt{T})$ regret bound for *online mirror*
44 *descent* (OMD), the latter does not directly yield regret bounds in the online setting. Moreover,
45 Orabona and Pál [2018] showed that OMD is not suitable when we do not know a priori the number
46 of iterations since it may suffer linear regret in this case. Finally, at present it is not known how
47 foundational OCO algorithms such as *follow the regularized leader* (FTRL) [Shalev-Shwartz, 2011,
48 Hazan, 2016] and *regularized dual averaging* [Xiao, 2010] (RDA) perform in the relative setting.

49 **Our results.** We analyze the performance of two general OCO algorithms: FTRL and dual-
50 stabilized OMD (DS-OMD, see [Fang et al., 2020]). We give $O(\sqrt{T})$ regret bounds in T rounds
51 for relative Lipschitz loss functions. Moreover, this is the first paper to show $O(\log T)$ regret if the
52 loss functions are also relative strongly-convex.¹ In addition, we are able to extend these bounds for
53 problems with composite loss functions, such as adding the ℓ_1 -norm to induce sparsity. The generality
54 of these algorithms lead to regret bounds for a wide variety of OCO algorithms (see Shalev-Shwartz
55 [2011], Hazan [2016] for some reductions). We demonstrate this flexibility by deriving convergence
56 rates for *dual averaging* Nesterov [2009] and *regularized dual averaging* [Xiao, 2010].

57 1.1 Related Work

58 Analyses of gradient descent methods in the differentiable convex setting usually require the objective
59 function f to be Lipschitz smooth, that is, the gradient of the objective function f is Lipschitz
60 continuous. Bauschke et al. [2017] proposed a generalized Lipschitz smoothness condition, called
61 *relative Lipschitz smoothness*, using Bregman divergences of a fixed reference function. They
62 proposed a proximal mirror descent method² called NoLips with a $O(1/T)$ convergence rate for
63 such functions. Building upon this work, Lu et al. [2018] slightly relaxed the definition of relative
64 smoothness and gave simpler analyses for mirror descent and dual averaging. Hanzely and Richtárik
65 [2018] propose and analyse coordinate and stochastic gradient descent methods for relatively smooth
66 functions. These ideas were later applied to non-convex problems by Mukkamala and Ochs [2019].
67 Unlike those prior works, in this paper we focus on the online case with non-differentiable loss
68 functions.

69 For non-differentiable convex optimization, Lipschitz continuity of the objective function is usually
70 needed to obtain a $O(1/\sqrt{T})$ convergence guarantee for classical methods. Lu [2019] showed that
71 this condition can be relaxed to what they called relative Lipschitz continuity of the objective function.
72 Under this latter assumption, they gave $O(1/\sqrt{T})$ convergence rates for deterministic and stochastic
73 mirror descent. In a similar vein, Grimmer [2019] showed how projected subgradient descent enjoys
74 a $O(1/\sqrt{T})$ convergence rate without Lipschitz continuity given that one has some control on the
75 norm of the subgradients. None of these works considered online algorithms. Although the results
76 from Lu [2019] for mirror descent can be adapted to the online setting, it is not clear how other
77 foundational OCO algorithms such as FTRL or RDA perform in this setting.

78 Antonakopoulos et al. [2020] generalized the Lipschitz continuity condition from the perspective of
79 Riemannian geometry. They proposed the notion of Riemannian Lipschitz-continuity (RLC) and
80 analyzed how OCO algorithms perform in this setting. They showed $O(\sqrt{T})$ regret bounds for both

¹This can be seen as analogous to the known logarithmic regret bounds when the loss functions are strongly convex [Hazan et al., 2007].

²They propose an algorithm in the general case with composite functions, but when we set $f := 0$ in their algorithm it boils down to classical mirror descent. In this case the novelty comes from the convergence analysis at a $O(1/T)$ rate without the use of classical Lipschitz smoothness.

81 FTRL and OMD with RLC cost functions in both the online and stochastic settings. However, the
 82 relationship between RLC and the relative Lipschitz-continuity notion of Lu [2019] is not clear.

83 Moreover, in the presence of both Lipschitz continuity and strong convexity we can obtain $O(1/T)$
 84 convergence rates in classical convex optimization [Bubeck, 2015, Section 3.4.1] and $O(\log T)$ regret
 85 in the online case [Hazan et al., 2007]. Lu [2019] extended the offline setting results by showing
 86 a $O(1/T)$ convergence rate when the objective function is both relative Lipschitz continuous and
 87 relative strongly convex. Still, this latter work does not obtain regret bounds for the online case. To
 88 the best of our knowledge, this is the first work studying conditions beyond strong convexity (and
 89 exp-concavity [Hazan et al., 2007]) to obtain logarithmic regret bounds.

90 2 Formal Definitions

91 Throughout this paper, \mathbb{R}^n denotes a n -dimensional real vector space endowed with an inner-product
 92 $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|$. We take $\mathcal{X} \subseteq \mathbb{R}^n$ to be a fixed closed convex set. The **dual norm** of $\|\cdot\|$ is defined
 93 by $\|x\|_* := \sup_{y \in \mathbb{R}^n: \|y\| \leq 1} \langle x, y \rangle$ for each $x \in \mathbb{R}^n$. Moreover, for any convex function $f: \mathcal{X} \rightarrow \mathbb{R}$
 94 and any $x \in \mathbb{R}^n$, a vector $g \in \mathbb{R}^n$ is a **subgradient** of f at x if g satisfies the *subgradient inequality*

$$f(z) \geq f(x) + \langle g, x - z \rangle, \quad \forall z \in \mathbb{R}^n. \quad (2.1)$$

95 We denote by $\partial f(x)$ the set of all subgradients of f at x , called the **subdifferential** of f at x . The
 96 **normal cone** of \mathcal{X} at a point $x \in \mathcal{X}$ is the set $N_{\mathcal{X}}(x) := \{a \in \mathbb{R}^n : \langle a, z - x \rangle \leq 0 \text{ for all } z \in \mathcal{X}\}$.

97 Let $R: \mathcal{D} \rightarrow \mathbb{R}$ be a convex function such that it is differentiable in $\mathcal{D}^\circ := \text{int } \mathcal{D}$ and such that we
 98 have $\mathcal{X} \subseteq \mathcal{D}^\circ$. The **Bregman divergence** (with respect to R) is given by

$$D_R(x, y) := R(x) - R(y) - \langle \nabla R(y), x - y \rangle, \quad \forall x \in \mathcal{D}, y \in \mathcal{D}^\circ.$$

99 An interesting and useful identity regarding Bregman divergences, sometimes called *three-point*
 100 *identity* [Bubeck, 2015], is

$$D_R(x, y) + D_R(z, x) - D_R(z, y) = \langle \nabla R(x) - \nabla R(y), x - z \rangle, \quad \forall z \in \mathcal{D}, \forall x, y \in \mathcal{D}^\circ. \quad (2.2)$$

101 Although the Bregman divergence with respect to R is not a metric, we can still interpret D_R as a
 102 way of measuring distances through the lens of R . An instructive example is the Bregman divergence
 103 associated with the squared ℓ_2 -norm $R := \frac{1}{2} \|\cdot\|_2^2$. In this case, we have $D_R(x, y) = \frac{1}{2} \|x - y\|_2^2$ for
 104 all $x, y \in \mathbb{R}^n$, that is, the divergence boils down to the squared ℓ_2 -distance. In light of this, a possible
 105 way to generalize Lipschitz continuity and strong convexity is to replace the norm in the classical
 106 definitions by the square root of the Bregman divergence [Lu et al., 2018].

107 First, recall that a function $f: \mathcal{X} \rightarrow \mathbb{R}$ is **L -Lipschitz continuous** with respect to $\|\cdot\|$ on $\mathcal{X}' \subseteq \mathcal{X}$ if

$$|f(x) - f(y)| \leq L \|x - y\|, \quad \forall x, y \in \mathcal{X}'.$$

108 Additionally, if f is convex, then the above definition implies³ that $\|g\|_* \leq L$ for all $x \in \mathcal{X}$ and all
 109 $g \in \partial f(x)$. Recall as well that a convex function $f: \mathcal{X} \rightarrow \mathbb{R}$ is **M -strongly convex** with respect to
 110 $\|\cdot\|$ on $\mathcal{X}' \subseteq \mathcal{X}$ for some $M > 0$ if

$$f(y) \geq f(x) + \langle g, y - x \rangle + \frac{M}{2} \|y - x\|^2, \quad \forall x, y \in \mathcal{X}', \forall g \in \partial f(x).$$

111 Let us now state generalizations of the above definitions due to Lu et al. [2018] and Lu [2019].

112 **Definition 2.1** (Relative Lipschitz continuity). A convex function $f: \mathcal{X} \rightarrow \mathbb{R}$ is **L -Lipschitz contin-**
 113 **uous** relative to R if

$$f(x) - f(y) \stackrel{(2.1)}{\leq} \langle g, x - y \rangle \leq L \sqrt{2D_R(y, x)}, \quad \forall x, y \in \mathcal{X}, \forall g \in \partial f(x). \quad (2.3)$$

114 The original definition of Lu [2019] requires $\|g\|_* \|x - y\| \leq L \sqrt{2D_R(x, y)}$ for all $x, y \in \mathcal{X}$ and
 115 $g \in \partial f(x)$. Since $\langle a, b \rangle \leq \|a\|_* \|b\|$ for any $a, b \in \mathbb{R}^n$, the above definition is slightly more general
 116 and does not depend on the choice of a norm.

³On the boundary of \mathcal{X} this implication is not as strong: we can only guarantee the existence of one subgradient with small norm. For our purposes this will not be of fundamental importance. For a more precise statement see [Ben-Tal and Nemirovski, 2001, §5.3]

117 **Definition 2.2** (Relative strong convexity [Lu et al., 2018]). A convex function $f: \mathcal{X} \rightarrow \mathbb{R}$ is
 118 **M -strongly convex** relative to R if

$$f(y) \geq f(x) + \langle g, y - x \rangle + MD_R(y, x), \quad \forall y, x \in \mathcal{X}, \forall g \in \partial f(x). \quad (2.4)$$

119 A notable special case of relative Lipschitz-continuity or relative strong convexity is when we pick
 120 $R := \frac{1}{2} \|\cdot\|_2^2$ and the classical definitions with respect to the ℓ_2 -norm are recovered.

121 **Example** (A function that is relative Lipschitz but not Lipschitz). Consider the function f given
 122 by $f(x) := x^2$ for each $x \in \mathbb{R}$. Since the derivative of f is unbounded on \mathbb{R} , it is not Lipschitz
 123 continuous on the entire line. Define the function R by $R(x) := 2x^4$ for all $x \in \mathbb{R}$. Then,

$$D_R(y, x) = 2y^4 - 2x^4 - 8x^3(y - x) = \frac{1}{2}(x^2 - y^2)^2 + x^2(x - y)^2 \geq x^2(x - y)^2, \quad \forall x, y \in \mathbb{R}.$$

124 Thus, $(f'(x)(x - y))^2 = 4x^2(x - y)^2 \leq 2 \cdot 2D_R(y, x)$ for any $x, y \in \mathbb{R}^n$. That is, f is $\sqrt{2}$ -Lipschitz
 125 continuous relative to R .

126 Lu [2019] discusses more substantial examples in detail, such as training of support vector machines,
 127 and finding a point in the intersection of several ellipsoids. Furthermore, he also gives a systematic
 128 way of picking a reference function for any objective functions whose subgradients at x have ℓ_2 -norm
 129 bounded by a polynomial in $\|x\|_2$. This useful construction allows many optimization problems to
 130 benefit from algorithms that are designed for the relative setting.

131 2.1 Conventions and Assumptions used Throughout the Paper

132 We collect here some additional notation and assumptions used throughout the paper.⁴ First, $\mathcal{X} \subseteq \mathbb{R}^n$
 133 denotes a closed convex set and $\{f_t\}_{t \geq 1}$ denotes a sequence of convex functions such that $f_t: \mathcal{X} \rightarrow \mathbb{R}$
 134 is subdifferentiable⁵ on \mathcal{X} for each $t \geq 1$. We denote by $\{\eta_t\}_{t \geq 0}$ a sequence of scalars such that
 135 $\eta_t \geq \eta_{t+1} > 0$ for each $t \geq 0$. Moreover, $\mathcal{D} \subseteq \mathbb{R}^n$ denotes a convex set with non-empty interior
 136 $\mathcal{D}^\circ := \text{int}(\mathcal{D})$ such that $\mathcal{X} \subseteq \mathcal{D}^\circ$. This latter set will be the domain of the regularizer of FTRL and
 137 of the mirror map for OMD. Namely, in Section 3 we denote by $R: \mathcal{D} \rightarrow \mathbb{R}$ the *regularizer* of FTRL,
 138 a convex function which is differentiable on \mathcal{D}° . In Section 5 we denote by $\Phi: \mathcal{D} \rightarrow \mathbb{R}$ the *mirror*
 139 *map* of online mirror descent (whose precise definition we postpone to Section 5).

140 3 Follow the Regularized Leader

141 *The follow the regularized leader* (FTRL) algorithm is a classical method for OCO. At each round,
 142 FTRL picks a point that minimizes the cost incurred by the previously seen functions plus a regularizer
 143 convex function (an *FTRL regularizer*). Intuitively, the latter helps the choices of the algorithm not
 144 to change too widely from one round to the next. In Algorithm 1 we formally outline the FTRL
 145 algorithm. It is well known [Hazan, 2016] that, in a game with T rounds, FTRL with properly tuned
 146 step sizes suffers at most $O(\sqrt{T})$ regret against Lipschitz continuous functions.⁶ When the loss
 147 functions are additionally strongly convex, FTRL suffers at most regret $O(\log T)$. In this section we
 148 describe one of our main results: the FTRL algorithm preserves these asymptotic regret guarantees in
 149 the relative setting.

150 The usual first step in the analyses of FTRL algorithms is to use basic properties of the iterates
 151 (without relying on convexity) to bound the algorithm’s regret by easier-to-analyse terms. Such
 152 bounds are usually the sum of two terms: the “diameter” of the feasible set through the lens of the
 153 FTRL regularizer and a sum of the difference in “quality” between consecutive iterates. For a classic
 154 example, see [Shalev-Shwartz, 2011, Lemma 2.3]. For our analysis we shall use a slightly tighter
 155 bound given by the Strong FTRL Lemma due to McMahan [2017]. For the sake of completeness we
 156 give a proof of this lemma (and discuss its applications in the composite setting) in Appendix B.1.

⁴The only exception is Lemma 3.1, which does not need convexity or differentiability of any of the functions.

⁵This is not too restrictive since convex functions are subdifferentiable on the relative interior of their domains [Rockafellar, 1997, Theorem 23.4].

⁶The big-O notation in this case hides constants that may depend on the dimension and other properties of the problem at hand. The best dependence on the Lipschitz constant and “distance to the comparison point” is usually achieved when the loss functions are Lipschitz continuous and the FTRL regularizer is strongly convex, both with respect to the same norm.

Algorithm 1 Follow the Regularized Leader (FTRL) Algorithm

Compute $x_1 \in \arg \min_{x \in \mathcal{X}} R(x)$
Set $F_0 := 0$
for $t = 1, 2, \dots$ **do**
 Observe f_t and suffer cost $f_t(x_t)$
 Set $F_t := F_{t-1} + f_t = \sum_{i=1}^t f_i$
 Compute $x_{t+1} \in \arg \min_{x \in \mathcal{X}} (F_t(x) + \frac{1}{\eta_t} R(x))$

157 **Lemma 3.1.** (Strong FTRL Lemma [McMahan, 2017]) Let $\{f_t\}_{t \geq 1}$ be a sequence of functions such
158 that $f_t: \mathcal{X} \rightarrow \mathbb{R}$ for each $t \geq 1$. Let $\{\eta_t\}_{t \geq 1}$ be a positive non-increasing sequence. Let $R: \mathcal{X} \rightarrow \mathbb{R}$
159 be such that $\{x_t\}_{t \geq 1}$ given as in Algorithm 1 is properly defined. If $F_t: \mathcal{X} \rightarrow \mathbb{R}$ is defined as in
160 Algorithm 1 for each $t \geq 1$, then,

$$\text{Regret}_T(z) \leq \sum_{t=0}^T \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) (R(z) - R(x_t)) + \sum_{t=1}^T (H_t(x_t) - H_t(x_{t+1})) \quad \forall T > 0,$$

161 where $\eta_0 := 1$, $\frac{1}{\eta_{-1}} := 0$, $x_0 := x_1$, and $H_t := F_t + \frac{1}{\eta_t} R$ for each $t \geq 1$.

162 3.1 Sublinear Regret with Relative Lipschitz Functions

163 In the following theorem we formally state our sublinear $O(\sqrt{T})$ regret bound of FTRL in T rounds
164 in the setting where the cost functions are Lipschitz continuous relative to the regularizer function
165 used in the FTRL method. The proof, which we defer to Appendix B.2, boils down to bounding
166 the terms $H_t(x_t) - H_t(x_{t+1})$ from the Strong FTRL Lemma by (roughly) $L^2 \eta_{t-1} / 2$. We do so
167 by combining the optimality conditions from the definition of the iterates in Algorithm 1 with the
168 L -Lipschitz continuity relative to R of the loss functions.

169 **Theorem 3.2.** Let $\{x_t\}_{t \geq 1}$ be defined as in Algorithm 1 and suppose f_t is L -Lipschitz continuous
170 relative to R for all $t \geq 1$. Let $z \in \mathcal{X}$ and let $K \in \mathbb{R}$ be such that $K \geq R(z) - R(x_1)$. Then,

$$\text{Regret}_T(z) \leq \frac{K}{\eta_T} + \sum_{t=1}^T \frac{L^2 \eta_{t-1}}{2}, \quad \forall T > 0.$$

171 In particular, if $\eta_t := \sqrt{K} / (L\sqrt{t+1})$ for each $t \geq 0$, then $\text{Regret}_T(z) \leq 2L\sqrt{K(T+1)}$.

172 3.2 Logarithmic Regret with Relative Strongly Convex Functions

173 Hazan et al. [2007] showed that if the cost functions are not only Lipschitz continuous but strongly
174 convex as well, then the *follow the leader* (FTL) method—FTRL without any regularizer—attains
175 logarithmic regret. Similarly, in this section we show that if the cost functions are relative Lipschitz
176 continuous and relative strongly convex, both relative to the same fixed function, then FTL suffers
177 regret at most logarithmic in the number of rounds. The proof of the next theorem is similar to the
178 proof of Theorem 3.2 and is deferred to Appendix B.3.

179 **Theorem 3.3.** Let $\{x_t\}_{t \geq 1}$ be defined as in Algorithm 1 with $R := 0$. Assume that f_t is L -Lipschitz
180 continuous and M -strongly convex relative to a differentiable convex function $h: \mathcal{D} \rightarrow \mathbb{R}$ for each
181 $t \geq 1$. Then, for all $z \in \mathcal{X}$,

$$\text{Regret}_T(z) \leq \frac{L^2}{2M} (\log(T) + 1), \quad \forall T > 0.$$

182 4 Dual Averaging and Composite Loss Functions

183 FTRL is a cornerstone algorithm in OCO, but sometimes it is not practical. Each iterate requires *exact*
184 minimization of the loss functions (plus the regularizer) which might not have always a closed form
185 solution. A notable special case of FTRL that mitigates this problem is the (online) *dual averaging*
186 (DA) method whose offline version is due to Nesterov [2009]. In each iteration, DA picks a point from
187 \mathcal{X} that minimizes the sum of past subgradients (scaled by the step size) plus a FTRL regularizer R .

188 Formally, for real convex functions $\{f_t\}_{t \geq 1}$ on \mathcal{X} , the online DA method computes iterates $\{x_t\}_{t \geq 1}$
 189 such that

$$x_{t+1} \in \arg \min_{x \in \mathcal{X}} \left(\eta_t \sum_{i=1}^t \langle g_i, x \rangle + R(x) \right) \quad \forall t \geq 0, \quad (4.1)$$

190 where $g_t \in \partial f_t(x_t)$ for each $t \geq 1$.

191 **Intuition.** It is well-known that the DA algorithm reduces to FTRL applied to the linearized
 192 functions $\{\tilde{f}_t\}_{t \geq 1}$ given by $\tilde{f}_t := \langle g_t, \cdot \rangle$ for each $t \in \mathbb{N}$ (for details see Hazan [2016, Lemma 5.4]).
 193 This reduction obviously preserves the property of being Lipschitz continuous since the gradient
 194 of \tilde{f}_t is g_t everywhere. A natural idea would be to use this same reduction in the relative setting.
 195 Unfortunately, this reduction does not preserve the property of being relative Lipschitz! Luckily,
 196 our proof only requires a weaker condition: being “relative Lipschitz” at the particular point x_t .
 197 Namely, the relative L -Lipschitzness (see (2.3)) of f_t implies $\langle \nabla \tilde{f}_t(x_t), x_t - y \rangle = \langle g_t, x_t - y \rangle \leq$
 198 $L\sqrt{2D_R(y, x_t)}$ for all $y \in \mathcal{X}$. That is all we need for the proof of Theorem 3.2 to go through,
 199 although we did state the theorem with this exact condition for the sake of simplicity. This discussion
 200 leads to the following corollary of Theorem 3.2.

201 **Corollary 4.1.** Let $\{x_t\}_{t \geq 1}$ be defined as in (4.1) and suppose f_t is L -Lipschitz continuous relative to
 202 R for all $t \geq 1$. Let $z \in \mathcal{X}$ and let $K \in \mathbb{R}$ be such that $K \geq R(z) - R(x_1)$. If $\eta_t := \sqrt{2K}/(L\sqrt{t+1})$
 203 for all $t \geq 1$, then $\text{Regret}_T(z) \leq 2L\sqrt{K(T+1)}$.

204 Another important consideration for applications is a variant of OCO in which the loss functions are
 205 composite [Duchi et al., 2010, Xiao, 2010]. More specifically, in this case we have a known “extra
 206 regularizer” Ψ , a (not necessarily differentiable) convex function, and add it to the loss functions. The
 207 goal is to induce some kind of structure in the iterates, such as adding ℓ_1 -regularization to promote
 208 sparsity. Note that OCO algorithms would still apply in this setting by replacing the loss functions f_t
 209 with $f_t + \Psi$ at each round t . However, in this case we are not exploiting the fact that the function Ψ is
 210 *known*. In the case of the relative setting, for example, it may be the case that the loss functions f_t are
 211 relative Lipschitz-continuous with respect to a certain function R , while Ψ is not. In Appendix C we
 212 extend the sublinear (composite) regret bound of Theorem 3.2 and show how this yields convergence
 213 bounds for regularized dual averaging [Xiao, 2010] in the relative setting.

214 5 Dual-Stabilized Online Mirror Descent

215 The mirror descent algorithm is a generalization of the classical gradient descent method that was
 216 first proposed by Nemirovsky and Yudin [1983]. A modern treatment was first given by Beck and
 217 Teboulle [2003]. The algorithm fits almost seamlessly into the OCO setting via a variant known
 218 as online mirror descent (OMD) (see [Hazan, 2016]). Recently, Orabona and Pál [2018] showed
 219 that OMD with a dynamic learning rate may suffer *linear* regret. (A dynamic learning rate is useful
 220 when we do not know the number of iterations ahead of time.) Moreover, this can happen even in
 221 simple and well-studied scenarios such as in the problem of prediction with expert advice, which
 222 corresponds to OMD equipped with negative entropy as a mirror map. In general, they showed that
 223 this may happen in cases where the Bregman divergence (with respect to the mirror map chosen) *is*
 224 *not* bounded over the entire feasible set. To resolve this issue, Fang et al. [2020] proposed a modified
 225 version of OMD called *dual-stabilized online mirror descent* (DS-OMD). In contrast to classical
 226 OMD, the regret bounds for the dual-stabilized version depend only on the Bregman divergence
 227 between the feasible set and the *initial iterate*.

228 We formally describe the DS-OMD method in Algorithm 2. Compared to OMD, DS-OMD adds an
 229 extra step in the dual space to mix the current dual iterate with the dual of the initial point. This step
 230 at iteration t is controlled by a stabilization parameter γ_t and it can be seen as a way to “stabilize” the
 231 algorithm in the dual space. Throughout this section we closely follow the notation and assumptions
 232 of Bubeck [2015, Chapter 4]. We assume that we have a **mirror map** for \mathcal{X} , that is, a differentiable
 233 strictly-convex function $\Phi: \mathcal{D} \rightarrow \mathbb{R}$ for \mathcal{X} such that the gradient of Φ diverges on the boundary of
 234 \mathcal{D} , that is, $\lim_{x \rightarrow \partial \mathcal{D}} \|\nabla \Phi(x)\|_2 = \infty$ where $\partial \mathcal{D} := \mathcal{D} \setminus \mathcal{D}^\circ$. These conditions on the mirror map
 235 guarantee that the algorithm is well-defined (for example, they guarantee the existence and uniqueness
 236 of the last step of Algorithm 2).

Algorithm 2 Dual-Stabilized Online Mirror Descent

Input: Stabilization coefficient γ_t and an initial iterate $x_1 \in \mathcal{X}$.

for $t = 1, 2, \dots$ **do**

 Observe f_t and suffer cost $f_t(x_t)$

 Compute $g_t \in \partial f_t(x_t)$

$\hat{x}_t := \nabla \Phi(x_t)$

$\hat{w}_{t+1} := \hat{x}_t - \eta_t g_t$

$\hat{y}_{t+1} := \gamma_t \hat{w}_{t+1} + (1 - \gamma_t) \hat{x}_1$

$y_{t+1} := \nabla \Phi^*(\hat{y}_{t+1})$

 Compute $x_{t+1} \in \arg \min_{x \in \mathcal{X}} D_\Phi(x, y_{t+1}) = \Phi(x) - \Phi(y_{t+1}) - \langle \nabla \Phi(y_{t+1}), x - y_{t+1} \rangle$

237 5.1 Sublinear Regret with Relative Lipschitz Functions

238 In this section, we give a regret bound for DS-OMD when the cost functions are all Lipschitz
239 continuous relative to the mirror map Φ . In this setting, if we set the stabilization coefficients to be
240 $\gamma_t := \eta_{t+1}/\eta_t$ and step size $O(1/\sqrt{t})$, DS-OMD obtains sublinear regret. This is formally stated in
241 the following theorem.

242 **Theorem 5.1.** Let $\{x_t\}_{t \geq 1}$ be defined as in Algorithm 2 with $\gamma_t := \eta_{t+1}/\eta_t$ for each $t \geq 1$. Assume
243 that f_t is L -Lipschitz continuous relative to Φ for all $t \geq 1$. Let $z \in \mathcal{X}$ and $K \in \mathbb{R}$ be such that
244 $K \geq D_\Phi(z, x_1)$. Then,

$$\text{Regret}_T(z) \leq \frac{K}{\eta_{T+1}} + \sum_{t=1}^T \frac{\eta_t L^2}{2}, \quad \forall T > 0.$$

245 In particular, if $\eta_t := \sqrt{K}/L\sqrt{t}$ for each $t \geq 1$, then $\text{Regret}_T(z) \leq 2L\sqrt{K(T+1)}$.

246 The proof is based on Theorem D.3, which gives an abstract regret upper bound for DS-OMD. Next
247 we compute specific upper bounds of $D_\Phi(x_t, w_{t+1})$ for each $t \geq 1$ by relative Lipschitz continuity
248 to make the abstract regret bound more specific. The whole proof of Theorem 5.1 is given in
249 Appendix D.1.

250 If we set each f_t to be a fixed function f and take average of all iterates, then we get the following
251 convergence rate for classical convex optimization as a corollary.

252 **Corollary 5.2.** Let Φ be a mirror map for \mathcal{X} and let $f: \mathcal{X} \rightarrow \mathbb{R}$ be a convex L -Lipschitz-continuous
253 function relative to Φ . Let $\{x_t\}_{t \geq 1}$ be given as in Algorithm 2 with loss functions $f_t := f$, step sizes
254 $\eta_t := \sqrt{K}/L\sqrt{t}$ for some $K \geq \sup_{z \in \mathcal{X}} D_\Phi(z, x_1)$, and stabilization parameter $\gamma_t := \eta_{t+1}/\eta_t$. If
255 $x^* \in \mathcal{X}$ is a minimizer of f , then,

$$f\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - f(x^*) \leq \frac{2L\sqrt{2K}}{\sqrt{T}}.$$

256 This recovers the same bound up to constant $4\sqrt{2}/3$ in Theorem 4.3 in Lu [2019], if we take $k = T - 1$
257 and $t_i = \frac{\sqrt{K}}{\sqrt{T}L}$ for $i \geq 0$ therein.

258 5.2 Logarithmic Regret with Relative Strongly Convex Functions

259 In Section 3.2 we showed that FTRL suffers at most logarithmic regret when the loss functions
260 are Lipschitz continuous and strongly convex, both relative to the same fixed reference function.
261 Similarly, we show that OMD suffers at most logarithmic regret if we have Lipschitz continuity and
262 strong convexity, both relative to the mirror map Φ . Interestingly, in this case the dual-stabilization
263 step can be skipped (that is, we can use $\gamma_t := 1$ for all t) and Algorithm 2 boils down to classic OMD.

264 **Theorem 5.3.** Let $\{x_t\}_{t \geq 1}$ be given as in Algorithm 2 with $\gamma_t := 1$ for all $t \geq 1$. Assume that f_t is
265 L -Lipschitz continuous and M -strongly convex relative to Φ for all $t \geq 1$. If $z \in \mathcal{X}$ and $\eta_t = \frac{1}{tM}$ for
266 each $t \geq 1$, then,

$$\text{Regret}_T(z) \leq \frac{L^2}{2M} (\log T + 1), \quad \forall T > 0.$$

267 The proof involves modifications of Theorem 5.1 and is deferred to Appendix D.2.

268 5.3 Sublinear Regret with Composite Loss Functions

269 We can extend our regret bounds to the setting with composite cost functions with minor modifications
270 to Algorithm 2. The classical version OMD adapted to this setting is due to Duchi et al. [2010] and is
271 known by composite objective mirror descent (COMID). They showed that COMID generalizes much
272 prior work like forward-backward splitting and derived new results on efficient matrix optimization
273 with Schatten p -norms based on this framework. Details of the modification needed on Algorithm D.3
274 in this setting together with regret bounds can be found in Appendix D.3.

275 6 Conclusions and Discussion

276 In this paper we showed regret bounds for both FTRL and stabilized OMD in the relative setting
277 proposed by Lu [2019]. All the results hold in the *anytime setting* in which we do not know the
278 number of rounds/iterations beforehand. Additionally, we gave logarithmic regret bounds for both
279 algorithms when the functions are relatively strongly convex, analogous to the results known in the
280 classical setting. Finally, we extend our results to the setting of composite cost functions, which
281 is pervasive in practice. These results open up the possibility of a new range of applications for
282 OCO algorithms and may allow for new analysis for known problems with better dependence on the
283 instance’s parameters.

284 At the moment there are at least two interesting directions for future research. The first would be to
285 investigate the connections among the different notions of relative smoothness, Lipschitz continuity,
286 and strong convexity in the literature. Another is to investigate systematic ways of choosing a
287 regularizer/mirror map for any given optimization problem. The latter was already an interesting
288 questions before notions of relative Lipschitz continuity and strong convexity were proposed, but
289 these new ideas give more flexibility in the choice of a regularizer.

290 7 Statement of Broader Impact

291 In this paper we study the performance of online convex optimization algorithms when the functions
292 are not necessarily Lipschitz continuous, a requirement in classical regret bounds. This opens up
293 the range of applications for which we can use OCO with good guarantees and guides how such
294 parameters such as regularizers/mirror maps and step sizes should be chosen. It is our hope that
295 this aids practitioners to develop more efficient ways to optimize and train their current models.
296 Furthermore, we hope theoreticians to be inspired to delve deep into the setting of non-smooth
297 optimization beyond Lipschitz continuity. It not only opens up the range of applications, but sheds
298 light onto the fundamental conditions on the cost functions and regularizers/mirror maps needed for
299 OCO algorithms to have good guarantees. Due to the theoretical nature of this work, we do not see
300 potentially bad societal or ethical impacts.

301 References

- 302 K. Antonakopoulos, E. V. Belmega, and P. Mertikopoulos. Online and stochastic optimization beyond
303 lipschitz continuity: A riemannian approach. 2020.
- 304 H. H. Bauschke, J. Bolte, and M. Teboulle. A descent lemma beyond lipschitz gradient continuity:
305 first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–
306 348, 2017.
- 307 A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex
308 optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- 309 A. Ben-Tal and A. Nemirovski. *Lectures on modern convex optimization*. MPS/SIAM Series on
310 Optimization. Society for Industrial and Applied Mathematics (SIAM), 2001.
- 311 S. Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine
312 Learning*, 8(3-4):231–357, 2015.
- 313 J. C. Duchi, S. Shalev-Shwartz, Y. Singer, and A. Tewari. Composite objective mirror descent. In
314 *COLT 2010*, pages 14–26. Omnipress, 2010.

- 315 J. C. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic
316 optimization. *J. Mach. Learn. Res.*, 12:2121–2159, 2011.
- 317 H. Fang, N. J. A. Harvey, V. S. Portella, and M. P. Friedlander. Online mirror descent and dual
318 averaging: keeping pace in the dynamic case. 2020. URL [https://arxiv.org/abs/2006.
319 02585](https://arxiv.org/abs/2006.02585).
- 320 B. Grimmer. Convergence rates for deterministic and stochastic subgradient methods without lipschitz
321 continuity. *SIAM Journal on Optimization*, 29(2):1350–1365, 2019.
- 322 F. Hanzely and P. Richtárik. Fastest rates for stochastic mirror descent methods. 2018. URL
323 <http://arxiv.org/abs/1803.07374>.
- 324 E. Hazan. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2
325 (3-4):157–325, 2016. URL <http://ocobook.cs.princeton.edu/OCObok.pdf>.
- 326 E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization.
327 *Machine Learning*, 69(2-3):169–192, 2007.
- 328 H. Lu. “Relative continuity” for non-lipschitz nonsmooth convex optimization using stochastic (or
329 deterministic) mirror descent. *Informs Journal on Optimization*, pages 265–352, 2019.
- 330 H. Lu, R. M. Freund, and Y. Nesterov. Relatively smooth convex optimization by first-order methods,
331 and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- 332 H. B. McMahan. A survey of algorithms and analysis for adaptive online learning. *The Journal of
333 Machine Learning Research*, 18(1):3117–3166, 2017.
- 334 M. C. Mukkamala and P. Ochs. Beyond alternating updates for matrix factorization with inertial
335 bregman proximal gradient algorithms. In *Advances in Neural Information Processing Systems*,
336 pages 4268–4278, 2019.
- 337 A. S. Nemirovsky and D. B. Yudin. Problem complexity and method efficiency in optimization.
338 1983.
- 339 Y. Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical programming*,
340 120(1):221–259, 2009.
- 341 F. Orabona and D. Pál. Scale-free online learning. *Theoretical Computer Science*, 716:50–69, 2018.
- 342 R. T. Rockafellar. *Convex analysis*. Princeton Landmarks in Mathematics. Princeton University Press,
343 Princeton, NJ, 1997. ISBN 0-691-01586-4. Reprint of the 1970 original, Princeton Paperbacks.
- 344 S. Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends® in
345 Machine Learning*, 4(2):107–194, 2011.
- 346 L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal
347 of Machine Learning Research*, 11(Oct):2543–2596, 2010.
- 348 M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Machine
349 Learning, Proceedings of the Twentieth International Conference (ICML 2003)*, pages 928–936,
350 2003.