

---

# Stochastic and Adversarial Online Learning without Hyperparameters

---

**Ashok Cutkosky**  
Department of Computer Science  
Stanford University  
ashokc@cs.stanford.edu

**Kwabena Boahen**  
Department of Bioengineering  
Stanford University  
boahen@stanford.edu

## Abstract

Most online optimization algorithms focus on one of two things: performing well in adversarial settings by adapting to unknown data parameters (such as Lipschitz constants), typically achieving  $O(\sqrt{T})$  regret, or performing well in stochastic settings where they can leverage some structure in the losses (such as strong convexity), typically achieving  $O(\log(T))$  regret. Algorithms that focus on the former problem hitherto achieved  $O(\sqrt{T})$  in the stochastic setting rather than  $O(\log(T))$ . Here we introduce an online optimization algorithm that achieves  $O(\log^4(T))$  regret in a wide class of stochastic settings while gracefully degrading to the optimal  $O(\sqrt{T})$  regret in adversarial settings (up to logarithmic factors). Our algorithm does not require any prior knowledge about the data or tuning of parameters to achieve superior performance.

## 1 Extending Adversarial Algorithms to Stochastic Settings

The online convex optimization (OCO) paradigm [1, 2] can be used to model a large number of scenarios of interest, such as streaming problems, adversarial environments, or stochastic optimization. In brief, an OCO algorithm plays  $T$  rounds of a game in which on each round the algorithm outputs a vector  $w_t$  in some convex space  $W$ , and then receives a loss function  $\ell_t : W \rightarrow \mathbb{R}$  that is convex. The algorithm’s objective is to minimize *regret*, which is the total loss of all rounds relative to  $w^*$ , the minimizer of  $\sum_{t=1}^T \ell_t$  in  $W$ :

$$R_T(w^*) = \sum_{t=1}^T \ell_t(w_t) - \ell_t(w^*)$$

OCO algorithms typically either make as few as possible assumptions about the  $\ell_t$  while attempting to perform well (adversarial settings), or assume that the  $\ell_t$  have some particular structure that can be leveraged to perform much better (stochastic settings). For the adversarial setting, the minimax optimal regret is  $O(BL_{\max}\sqrt{T})$ , where  $B$  is the diameter of  $W$  and  $L_{\max}$  is the maximum Lipschitz constant of the losses [3]. A wide variety of algorithms achieve this bound without prior knowledge of one or both of  $B$  and  $L_{\max}$  [4, 5, 6, 7], resulting in hyperparameter-free algorithms. In the stochastic setting, it was recently shown that for a class of problems (those satisfying the so-called *Bernstein condition*), one can achieve regret  $O(dBL_{\max}\log(T))$  where  $W \subset \mathbb{R}^d$  using the METAGRAD algorithm [8, 9]. This approach requires knowledge of the parameter  $L_{\max}$ .

In this paper, we extend an algorithm for the parameter-free adversarial setting [7] to the stochastic setting, achieving both optimal regret in adversarial settings as well as logarithmic regret in a wide class of stochastic settings, without needing to tune parameters. Our class of stochastic settings is those for which  $\mathbb{E}[\nabla\ell_t(w_t)]$  is aligned with  $w_t - w^*$ , quantified by a value  $\alpha$  that increases with

increasing alignment. We call losses in this class  $\alpha$ -acutely convex, and show that a single quadratic lower bound on the average loss is sufficient to ensure high  $\alpha$ .

This paper is organized as follows. In Section 2, we provide an overview of our approach. In Section 3, we give explicit pseudo-code and prove our regret bounds for the adversarial setting. In Section 4, we formally define  $\alpha$ -acute convexity and prove regret bounds for the acutely convex stochastic setting. Finally, in Section 5, we give some motivating examples of acutely convex stochastic losses. Section 6 concludes the paper.

## 2 Overview of Approach

Before giving the overview, we fix some notation. We assume our domain  $W$  is a closed convex subset of a Hilbert space with  $0 \in W$ . We write  $g_t$  to be an arbitrary subgradient of  $\ell_t$  at  $w_t$  for all  $t$ , which we denote by  $g_t \in \partial\ell_t(w_t)$ .  $L_{\max}$  is the maximum Lipschitz constant of all the  $\ell_t$ , and  $B$  is the diameter of the space  $W$ . The norm  $\|\cdot\|$  we use is the 2-norm:  $\|w\| = \sqrt{w \cdot w}$ . We observe that since each  $\ell_t$  is convex, we have  $R_T(w^*) \leq \sum_{t=1}^T g_t(w_t - w^*)$ . We will make heavy use of this inequality; every regret bound we state will in fact be an upper bound on  $\sum_{t=1}^T g_t(w_t - w^*)$ . Finally, we use a compressed sum notation  $g_{1:t} = \sum_{t'=1}^t g_{t'}$ , and we use  $\tilde{O}$  to suppress logarithmic terms in big-Oh notation. All proofs omitted from the main text appear in the appendix.

Our algorithm works by trading off some performance in order to avoid knowledge of problem parameters. Prior analysis of the METAGRAD algorithm [9] showed that any algorithm guaranteeing  $R_T(w^*) = \tilde{O}\left(\sqrt{\sum_{t=1}^T (g_t \cdot (w_t - w^*))^2}\right)$  will obtain logarithmic regret for stochastic settings satisfying the Bernstein condition. We will instead guarantee the weaker regret bound:

$$R_T(w^*) \leq \tilde{O}\left(\sqrt{L_{\max} \sum_{t=1}^T \|g_t\| \|w_t - w^*\|^2}\right) \quad (1)$$

which we will show in turn implies  $\sqrt{T}$  regret in adversarial settings and logarithmic regret for acutely convex stochastic settings. Although (1) is weaker than the METAGRAD regret bound, we can obtain it without prior knowledge.

In order to come up with an algorithm that achieves the bound (1), we interpret it as the square root of  $\mathbb{E}[\|w - w^*\|^2]$ , where  $w$  takes on value  $w_t$  with probability proportional to  $\|g_t\|$ . This allows us to use the bias-variance decomposition to write (1) as:

$$R_T(w^*) \leq \tilde{O}\left(\|w^* - \bar{w}\| \sqrt{L_{\max} \|g\|_{1:T}} + \sqrt{\sum_{t=1}^T L_{\max} \|g_t\| \|w_t - \bar{w}\|^2}\right) \quad (2)$$

where  $\bar{w} = \frac{\sum_{t=1}^T \|g_t\| w_t}{\|g\|_{1:T}}$ . Certain algorithms for unconstrained OCO can achieve  $R_T(u) = \tilde{O}(\|u\| L_{\max} \sqrt{\|g\|_{1:T}})$  simultaneously for all  $u \in W$  [10, 6, 11, 7]. Thus if we knew  $\bar{w}$  ahead of time, we could translate the predictions of one such algorithm by  $\bar{w}$  to obtain  $R_T(w^*) \leq \tilde{O}(\|w^* - \bar{w}\| L_{\max} \sqrt{\|g\|_{1:T}})$ , the bias term of (2). We do not know  $\bar{w}$ , but we can estimate it over time. Errors in the estimation procedure will cause us to incur the variance term of (2). We implement this strategy by modifying FREEREX [7], an unconstrained OCO algorithm that does not require prior knowledge of any parameters.

Our modification to FREEREX is very simple: we set  $w_t = \hat{w}_t + \bar{w}_{t-1}$  where  $\hat{w}_t$  is the  $t^{\text{th}}$  output of FREEREX, and  $\bar{w}_{t-1}$  is (approximately) a weighted average of the previous vectors  $w_1, \dots, w_{t-1}$  with the weight of  $w_t$  equal to  $\|g_t\|$ . This  $\bar{w}_t$  offset can be viewed as a kind of momentum term that accelerates us towards optimal points when the losses are stochastic (which tends to cause correlated  $w_t$  and therefore large offsets), but has very little effect when the losses are adversarial (which tends to cause uncorrelated  $w_t$  and therefore small offsets).

### 3 FREEREXMOMENTUM

In this section, we explicitly describe and analyze our algorithm, FREEREXMOMENTUM, a modification of FREEREX. FREEREX is a Follow-the-Regularized-Leader (FTRL) algorithm, which means that for all  $t$ , there is some regularizer function  $\psi_t$  such that  $w_{t+1} = \operatorname{argmin}_W \psi_t(w) + g_{1:t} \cdot w$ . Specifically, FREEREX uses  $\psi_t = \frac{\sqrt{5}}{a_t \eta_t} \phi(a_t w)$ , where  $\phi(w) = (\|w\| + 1) \log(\|w\| + 1) - \|w\|$  and  $\eta_t$  and  $a_t$  are specific numbers that grow over time as specified in Algorithm 1. FREEREXMOMENTUM's predictions are given by offsetting FREEREX's predictions  $w_{t+1}$  by a momentum term  $\bar{w}_t = \frac{\sum_{t'=1}^{t-1} \|g_{t'}\| w_{t'}}{1 + \|g\|_{1:t}}$ . We accomplish this by shifting the regularizers  $\psi_t$  by  $\bar{w}_t$ , so that FREEREXMOMENTUM is FTRL with regularizers  $\psi_t(w - \bar{w}_t)$ .

---

#### Algorithm 1 FREEREXMOMENTUM

---

**Initialize:**  $\frac{1}{\eta_0^2} \leftarrow 0$ ,  $a_0 \leftarrow 0$ ,  $w_1 \leftarrow 0$ ,  $L_0 \leftarrow 0$ ,  $\psi(w) = (\|w\| + 1) \log(\|w\| + 1) - \|w\|$   
**for**  $t = 1$  **to**  $T$  **do**  
  Play  $w_t$   
  Receive subgradient  $g_t \in \partial \ell_t(w_t)$   
   $L_t \leftarrow \max(L_{t-1}, \|g_t\|)$ . //  $L_t = \max_{t' \leq t} \|g_{t'}\|$   
   $\frac{1}{\eta_t^2} \leftarrow \max\left(\frac{1}{\eta_{t-1}^2} + 2\|g_t\|^2, L_t \|g_{1:t}\|\right)$ .  
   $a_t \leftarrow \max(a_{t-1}, 1/(L_t \eta_t)^2)$   
   $\bar{w}_t \leftarrow \frac{\sum_{t'=1}^{t-1} \|g_{t'}\| w_{t'}}{1 + \|g\|_{1:t}}$   
   $w_{t+1} \leftarrow \operatorname{argmin}_W \left[ \frac{\sqrt{5} \phi(a_t(w - \bar{w}_t))}{a_t \eta_t} + g_{1:t} \cdot w \right]$   
**end for**

---

#### 3.1 Regret Analysis

We leverage the description of FREEREXMOMENTUM in terms of shifted regularizers to prove a regret bound of the same form as (1) in four steps:

1. From [7] Theorem 13, we bound the regret by

$$\begin{aligned} R_T(w^*) &\leq \sum_{t=1}^T g_t \cdot (w_t - w^*) \\ &\leq \psi_T(w^*) + \sum_{t=1}^T \psi_{t-1}(w_{t+1}^+) - \psi_t^+(w_{t+1}^+) + g_t \cdot (w_t - w_{t+1}^+) \\ &\quad + \psi_T^+(w^*) - \psi_T(w^*) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) \end{aligned}$$

where  $\psi_t^+(w) \approx \frac{\sqrt{5} \phi(a_t(w - \bar{w}_{t-1}))}{a_t \eta_t}$  is a version of  $\psi_t$  shifted by  $\bar{w}_{t-1}$  instead of  $\bar{w}_t$ , and  $w_{t+1}^+ = \operatorname{argmin}_W \psi_t^+(w) + g_{1:t} w$ . This breaks the regret out into two sums, one in which we have the term  $\psi_{t-1}(w_{t+1}^+) - \psi_t^+(w_{t+1}^+)$  for which the two different functions are shifted by the same amount, and one with the term  $\psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+)$ , for which the functions are shifted differently, but the arguments are the same.

2. Because  $\psi_{t-1}$  and  $\psi_t^+$  are shifted by the same amount, the regret analysis for FREEREX in [7] applies to the second line of the regret bound, yielding a quantity similar to  $\|w^* - \bar{w}_T\| \sqrt{L_{\max} \|g\|_{1:T}}$ .
3. Next, we analyze the third line. We show that  $\bar{w}_t - \bar{w}_{t-1}$  cannot be too big, and use this observation to bound the third line with a quantity similar to  $\sqrt{\sum_{t=1}^T L_{\max} \|g_t\| (w_t - \bar{w}_T)^2}$ . At this point we have enough results to prove a bound of the form (2) (see Theorem 1).
4. Finally, we perform some algebraic manipulation on the bound from the first three steps to obtain a bound of the form (1) (see Corollary 2).

The details of Steps 1-3 procedure are in the appendix, resulting in Theorem 1, stated below. Step 4 is carried out in Corollary 2, which follows.

**Theorem 1.** Let  $\psi(w) = (\|w\|+1) \log(\|w\|+1) - \|w\|$ . Set  $L_t = \max_{t' \leq t} \|g_{t'}\|$ , and  $Q_T = 2 \frac{\|g\|_{1:T}}{L_{\max}}$ . Define  $\frac{1}{\eta_t}$  and  $a_t$  as in the pseudo-code for FREEREXMOMENTUM (Algorithm 1). Then the regret of FREEREXMOMENTUM is bounded by:

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq \frac{\sqrt{5}}{Q_T \eta_T} \psi(Q_T(w^* - \bar{w}_T)) + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(Ba_T + 1) \\ &+ \sqrt{2L_{\max} \left( \|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\| \|w_t - \bar{w}_T\|^2 \right)} \left( 2 + \log \left( \frac{1 + \|g\|_{1:T}}{1 + \|g_1\|} \right) \right) \log(Ba_T + 1) \end{aligned}$$

**Corollary 2.** Under the assumptions and notation of Theorem 1, the regret of FREEREXMOMENTUM is bounded by:

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right)} \log(2BT + 1) (2 + \log(T)) \\ &+ 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(2BT + 1) \end{aligned}$$

Observe that since  $w_t$  and  $w^*$  are both in  $W$ ,  $\|w^*\|$  and  $\|w_t - w^*\|$  both are at most  $B$ , so that Corollary 2 implies that FREEREXMOMENTUM achieves  $\tilde{O}(BL_{\max}\sqrt{T})$  regret in the worst-case, which is optimal up to logarithmic factors.

### 3.2 Efficient Implementation for $L_{\infty}$ Balls

A careful reader may notice that the procedure for FREEREXMOMENTUM involves computing  $\operatorname{argmin}_W \left[ \frac{\sqrt{5}\psi(a_t(w - \bar{w}_t))}{a_t \eta_t} + g_{1:t} \cdot w \right]$ , which may not be easy if the solution  $w_{t+1}$  is on the boundary of  $W$ . When the  $w_{t+1}$  is not on the boundary of  $W$ , then we have a closed-form update:

$$w_{t+1} = \bar{w}_t - \frac{g_{1:t}}{a_t \|g_{1:t}\|} \left[ \exp \left( \frac{\eta_t \|g_{1:t}\|}{\sqrt{5}} \right) - 1 \right] \quad (3)$$

However, when  $w_{t+1}$  lies on the boundary of  $W$ , it is not clear how to compute it for general  $W$ . In this section we offer a simple strategy for the case that  $W$  is an  $L_{\infty}$  ball,  $W = \prod_{i=1}^d [-b, b]$ .

In this setting, we can use the standard trick (e.g. see [12]) of running a separate copy of FREEREXMOMENTUM for each coordinate. That is, we observe that

$$R_T(w^*) \leq \sum_{t=1}^T g_t \cdot (w_t - u) = \sum_{i=1}^d \sum_{t=1}^T g_{t,i} (w_{t,i} - u_i) \quad (4)$$

so that if we run an independent online learning algorithm on each coordinate, using the coordinates of the gradients  $g_{t,i}$  as losses, then the total regret is at most the sum of the individual regrets. More detailed pseudocode is given in Algorithm 2.

Coordinate-wise FREEREXMOMENTUM is easily implementable in time  $O(d)$  per update because the FREEREXMOMENTUM update is easy to perform in one dimension: if the update (3) is outside the domain  $[-b, b]$ , simply set  $w_{t+1}$  to  $b$  or  $-b$ , whichever is closer to the unconstrained update. Therefore, coordinate-wise FREEREXMOMENTUM can be computed in  $O(d)$  time per update.

We bound the regret of coordinate-wise FREEREXMOMENTUM using Corollary 2 and Equation (4), resulting the following Corollary.

---

**Algorithm 2** Coordinate-Wise FREEREXMOMENTUM

---

**Initialize:**  $w_1 = 0$ ,  $d$  copies of FREEREXMOMENTUM,  $F_1, \dots, F_d$ , where each  $F_i$  uses domain  $W = [-b, b]$ .  
**for**  $t = 1$  **to**  $T$  **do**  
    Play  $w_t$ , receive subgradient  $g_t$ .  
    **for**  $i = 1$  **to**  $d$  **do**  
        Give  $g_{t,i}$  to  $F_i$ .  
        Get  $w_{t+1,i} \in [-b, b]$  from  $F_i$ .  
    **end for**  
**end for**

---

**Corollary 3.** *The regret of coordinate-wise FREEREXMOMENTUM is bounded by:*

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq 2\sqrt{5} \sqrt{dL_{\max} \left( d\|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right) \log(2Tb + 1)(2 + \log(T))} \\ &\quad + 405dL_{\max} + 2L_{\max}db + 3d \frac{L_{\max}\sqrt{2L_{\max}}}{\sqrt{1+L_1}} b \log(2bT + 1) \end{aligned}$$

## 4 Logarithmic Regret in Stochastic Problems

In this section we formally define  $\alpha$ -acute convexity and show that FREEREXMOMENTUM achieves logarithmic regret for  $\alpha$ -acutely convex losses. As a warm-up, we first consider the simplest case in which the loss functions  $\ell_t$  are fixed,  $\ell_t = \ell$  for all  $t$ . After showing logarithmic regret for this case, we will then generalize to more complicated stochastic settings.

Intuitively, an acutely convex loss function  $\ell$  is one for which the gradient  $g_t$  is aligned with the vector  $w_t - w^*$  where  $w^* = \operatorname{argmin} \ell$ , as defined below.

**Definition 4.** *A convex function  $\ell$  is  $\alpha$ -acutely convex on a set  $W$  if  $\ell$  has a global minimum at some  $w^* \in W$  and for all  $w \in W$ , for all subgradients  $g \in \partial\ell(w)$ , we have*

$$g \cdot (w - w^*) \geq \alpha \|g\| \|w - w^*\|^2$$

With this definition in hand, we can show logarithmic regret in the case where  $\ell_t = \ell$  for all  $t$  for some  $\alpha$ -acutely convex function  $\ell$ . From Corollary 2, with  $w^* = \operatorname{argmin} \ell$ , we have

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq \tilde{O} \left( \sqrt{L_{\max} \left( \|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right)} \right) \\ &\leq \tilde{O} \left( \sqrt{L_{\max} \left( \|w^*\| + \frac{1}{\alpha} \sum_{t=1}^T g_t \cdot (w^* - w_t) \right)} \right) \end{aligned} \quad (5)$$

Where the  $\tilde{O}$  notation suppresses terms whose dependence on  $T$  is at most  $O(\log^2(T))$ . Now we need a small Proposition:

**Proposition 5.** *If  $a, b, c$  and  $d$  are non-negative constants such that*

$$x \leq a\sqrt{bx + c} + d$$

*Then*

$$x \leq 4a^2b + 2a\sqrt{c} + 2d$$

Applying Proposition 5 to Equation (5) with  $x = \sum_{t=1}^T g_t \cdot (w_t - w^*)$  yields

$$R_T(u) \leq \tilde{O} \left( \frac{L_{\max}\|w^*\|}{\alpha} \right)$$

where the  $\tilde{O}$  again suppresses logarithmic terms, now with dependence on  $T$  at most  $O(\log^4(T))$ .

Having shown that FREEREXMOMENTUM achieves logarithmic regret on fixed  $\alpha$ -acutely convex losses, we now generalize to stochastic losses. In order to do this we will necessarily have to make some assumptions about the process generating the stochastic losses. We encapsulate these assumptions in a stochastic version of  $\alpha$ -acute convexity, given below.

**Definition 6.** Suppose for all  $t$ ,  $g_t$  is such that  $\mathbb{E}[g_t | g_1, \dots, g_{t-1}] \in \partial \ell(w_t)$  for some convex function  $\ell$  with minimum at  $w^*$ . Then we say  $g_t$  is  $\alpha$ -acutely convex in expectation if:

$$\mathbb{E}[g_t] \cdot (w_t - w^*) \geq \alpha \mathbb{E}[\|g_t\| \|w_t - w^*\|^2]$$

where all expectations are conditioned on  $g_1, \dots, g_{t-1}$ .

Using this definition, a fairly straightforward calculation gives us the following result.

**Theorem 7.** Suppose  $g_t$  is  $\alpha$ -acutely convex in expectation and  $g_t$  is bounded  $\|g_t\| \leq L_{\max}$  with probability 1. Then FREEREXMOMENTUM achieves expected regret:

$$\mathbb{E}[R_T(w^*)] \leq \tilde{O} \left( \frac{L_{\max} \|w^*\|}{\alpha} \right)$$

*Proof.* Throughout this proof, all expectations are conditioned on prior subgradients. By Corollary 2 and Jensen's inequality we have

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T g_t \cdot (w_t - w^*) \right] &\leq \mathbb{E} \left[ 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(2BT+1) \right. \\ &\quad \left. + 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right)} \log(2TB+1)(2+\log(T)) \right] \\ &\leq 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{\delta}} B \log(2BT+1) \\ &\quad + 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \sum_{t=1}^T \mathbb{E}[\|g_t\| \|w^* - w_t\|^2] \right)} \log(2TB+1)(2+\log(T)) \\ &\leq 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{\delta}} B \log(2BT+1) \\ &\quad + 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \frac{1}{\alpha} \sum_{t=1}^T \mathbb{E}[g_t \cdot (w_t - w^*)] \right)} \log(2TB+1)(2+\log(T)) \end{aligned}$$

Set  $R = \mathbb{E} \left[ \sum_{t=1}^T g_t(w_t - w^*) \right]$ . Then we have shown

$$\begin{aligned} R &\leq 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \frac{R}{\alpha} \right)} \log(2TB+1)(2+\log(T)) \\ &\quad + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{\delta}} B \log(BT+1) \\ &= \tilde{O} \left[ \sqrt{L_{\max} \left( \|w^*\|^2 + \frac{R}{\alpha} \right)} \right] \end{aligned}$$

And now we use Proposition 5 to conclude:

$$\sum_{t=1}^T \mathbb{E}[g_t \cdot (w_t - w^*)] = \tilde{O} \left( \frac{L_{\max} \|w^*\|}{\alpha} \right)$$

as desired, where again  $\tilde{O}$  hides at most a  $O(\log^4(T))$  dependence on  $T$ .  $\square$

Exactly the same argument with an extra factor of  $d$  applies to the regret of FREEREXMOMENTUM with coordinate-wise updates.

## 5 Examples of $\alpha$ -acute convexity in expectation

In this section, we show that  $\alpha$ -acute convexity in expectation is a condition that arises in practice, justifying the relevance of our logarithmic regret bounds. To do this, we show that a quadratic lower bound on the expected loss implies  $\alpha$ -acute convexity, demonstrating acute convexity is a weaker condition than strong convexity.

**Proposition 8.** *Suppose  $\mathbb{E}[g_t | g_1, \dots, g_{t-1}] \in \partial \ell(w_t)$  for some convex  $\ell$  such that for some  $\mu > 0$  and  $w^* = \operatorname{argmin} \ell$ ,  $\ell(w) - \ell(w^*) \geq \frac{\mu}{2} \|w - w^*\|^2$  for all  $w \in W$ . Suppose  $\|g\| \leq L_{\max}$  with probability 1. Then  $g_t$  is  $\frac{\mu}{2L_{\max}}$ -acutely convex in expectation.*

*Proof.* By convexity and the hypothesis of the proposition:  $\mathbb{E}[g_t] \cdot (w_t - w^*) \geq \ell(w_t) - \ell(w^*) \geq \frac{\mu}{2} \|w_t - w^*\|^2 \geq \frac{\mu}{2L_{\max}} \mathbb{E}[\|g_t\|] \|w_t - w^*\|^2$   $\square$

With Proposition 8, we see that FREEREXMOMENTUM obtains logarithmic regret for any loss that is larger than a quadratic, without requiring knowledge of the parameter  $\mu$  or the Lipschitz bound  $L_{\max}$ . Further, this result requires only the *expected loss*  $\ell = \mathbb{E}[\ell_t]$  to have a quadratic lower bound - the individual losses  $\ell_t$  themselves need not do so.

The boundedness of  $W$  makes it surprisingly easy to have a quadratic lower bound. Although a quadratic lower bound for a function  $\ell$  is easily implied by strong convexity, the quadratic lower bound is a significantly weaker condition. For example, since  $W$  has diameter  $B$ ,  $\|w\| \geq \frac{1}{B} \|w\|^2$  and so the absolute value is  $\frac{1}{B}$ -acutely convex, but not strongly convex. The following Proposition shows that existence of a quadratic lower bound is actually a *local* condition; so long as the expected loss  $\ell$  has a quadratic lower bound in a neighborhood of  $w^*$ , it must do so over the entire space  $W$ :

**Proposition 9.** *Suppose  $\ell : W \rightarrow \mathbb{R}$  is a convex function such that  $\ell(w) - \ell(w^*) \geq \frac{\mu}{2} \|w - w^*\|^2$  for all  $w$  with  $\|w - w^*\| \leq r$ . Then  $\ell(w) - \ell(w^*) \geq \min\left(\frac{\mu r}{2B}, \frac{\mu}{2}\right) \|w - w^*\|^2$  for all  $w \in W$ .*

*Proof.* We translate by  $w^*$  to assume without loss of generality that  $w^* = 0$ . Then the statement is clear for  $\|w\| \leq r$ . By convexity,  $\ell(w) - \ell(w^*) \geq \frac{\|w\|}{r} \left[ \ell\left(\frac{rw}{\|w\|}\right) - \ell(w^*) \right] \geq \frac{\mu r}{2} \|w\| \geq \frac{\mu r}{2B} \|w\|^2$ .  $\square$

Finally, we provide a simple motivating example of an interesting problem we can solve with an  $\alpha$ -acutely convex loss that is not strongly convex: computing the median.

**Proposition 10.** *Let  $W = [a, b]$ , and  $\ell_t(w) = |w - x_t|$  where each  $x_t$  is drawn i.i.d. from some fixed distribution with a continuous cumulative distribution function  $D$ , and assume  $D(x^*) = \frac{1}{2}$ . Further, suppose  $|2D(w) - 1| \geq F|w - x^*|$  for all  $|w - x^*| \leq G$ . Suppose  $g_t = \ell'_t(w_t)$  for  $w_t \neq x_t$  and  $g_t = \pm 1$  with equal probability if  $w_t = x_t$ . Then  $g_t$  is  $\min\left(\frac{FG}{b-a}, F\right)$ -acutely convex in expectation.*

*Proof.* By a little calculation,  $\mathbb{E}[g_t] = \ell'(w_t) = 2D(w_t) - 1$ , and  $\mathbb{E}[|g_t|] = 1$ . Since  $\ell'(x^*) = 0$ ,  $w^* = x^*$  (the median). For  $|w_t - x^*| \geq G$ , we have  $|2D(w) - 1| \geq FG$ , which gives  $\mathbb{E}[g_t] \cdot (w_t - w^*) \geq \frac{FG}{b-a} \mathbb{E}[|g_t|] (w_t - w^*)^2$ . For  $|w_t - x^*| \leq G$ , we have  $\mathbb{E}[g_t] \cdot (w_t - w^*) \geq F \mathbb{E}[|g_t|] (w_t - w^*)^2$ , so that  $g_t$  is  $\min\left(\frac{FG}{b-a}, F\right)$ -acutely convex in expectation.  $\square$

Proposition 10 shows that we can obtain low regret for an interesting stochastic problem without curvature. The condition on the cumulative distribution function  $D$  is asking only that there be positive density in a neighborhood of the median; it would be satisfied if  $D'(w) \geq F$  for  $|w| \leq G$ .

If the expected loss  $\ell$  is  $\mu$ -strongly convex, we can apply Proposition 8 to see that  $\ell$  is  $\mu/2$ -aligned, and then use Theorem 7 to obtain a regret of  $\tilde{O}(L_{\max} \|w^*\|/\mu)$ . This is different from the usual regret bound of  $\tilde{O}(L_{\max}^2/\mu)$  obtained by Online Newton Step [13], which is due to an inefficiency in using the weaker  $\alpha$ -alignment condition. Instead, arguing from the regret bound of Corollary 2 directly, we can recover the optimal regret bound:

**Corollary 11.** *Suppose each  $\ell_t$  is an independent random variable with  $\mathbb{E}[\ell_t] = \ell$  for some  $\mu$ -strongly convex  $\ell$  with minimum at  $w^*$ . Then the expected regret of FREEREXMOMENTUM satisfies*

$$\mathbb{E} \left[ \sum_{t=1}^T \ell(w_t) - \ell(w^*) \right] \leq \tilde{O}(L_{\max}^2/\mu)$$

Where the  $\tilde{O}$  hides terms that are logarithmic in  $TB$ .

*Proof.* From strong-convexity, we have

$$\|w_t - w^*\|^2 \leq \frac{2}{\mu}(\ell(w_t) - \ell(w^*))$$

Therefore applying Corollary 2 we have

$$\begin{aligned} \mathbb{E}[R_T(w^*)] &= \mathbb{E} \left[ \sum_{t=1}^T \ell(w_t) - \ell(w^*) \right] \leq \tilde{O} \left( \sqrt{L_{\max}^2 \mathbb{E} \left[ \sum_{t=1}^T \|w_t - w^*\|^2 \right]} \right) \\ &\leq \tilde{O}(\sqrt{L_{\max}^2 \mathbb{E}[R_T(w^*)]}) \end{aligned}$$

So that applying Proposition 5 we obtain the desired result.  $\square$

As a result of Corollary 11, we see that FREEREXMOMENTUM obtains logarithmic regret for  $\alpha$ -aligned problems and also obtains the optimal (up to log factors) regret bound for  $\mu$ -strongly-convex problems, all without requiring any knowledge of the parameters  $\alpha$  or  $\mu$ . This stands in contrast to prior algorithms that adapt to user-supplied curvature information such as Adaptive Gradient Descent [14] or  $(\mathcal{A}, \mathcal{B})$ -prod [15].

## 6 Conclusions and Open Problems

We have presented an algorithm, FREEREXMOMENTUM, that achieves both  $\tilde{O}(BL_{\max}\sqrt{T})$  regret in adversarial settings and  $\tilde{O}(\frac{L_{\max}B}{\alpha})$  regret in  $\alpha$ -acutely convex stochastic settings without requiring any prior information about any parameters. We further showed that a quadratic lower bound on the expected loss implies acute convexity, so that while strong-convexity is sufficient for acute convexity, other important loss families such as the absolute loss may also be acutely convex. Since FREEREXMOMENTUM does not require prior information about any problem parameters, it does not require any hyperparameter tuning to be assured of good convergence. Therefore, the user need not actually know whether a particular problem is adversarial or acutely convex and stochastic, or really much of anything at all about the problem, in order to use FREEREXMOMENTUM.

There are still many interesting open questions in this area. First, we would like to find an efficient way to implement the FREEREXMOMENTUM algorithm or some variant directly, without appealing to coordinate-wise updates. This would enable us to remove the factor of  $d$  we incur by using coordinate-wise updates. Second, our modification to FREEREX is extremely simple and intuitive, but our analysis makes use of some of the internal logic of FREEREX. It is possible, however, that *any* algorithm with sufficiently low regret can be modified in a similar way to achieve our results. Finally, we observe that while  $\log^4(T)$  is much better than  $\sqrt{T}$  asymptotically, it turns out that  $\log^4(T) > \sqrt{T}$  for  $T < 10^{11}$ , which casts the practical relevance of our logarithmic bounds in doubt. Therefore we hope that this work serves as a starting point for either new analysis or algorithm design that further simplifies and improves regret bounds.

## References

- [1] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 928–936, 2003.
- [2] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2011.



- [3] Jacob Abernethy, Peter L Bartlett, Alexander Rakhlin, and Ambuj Tewari. Optimal strategies and minimax lower bounds for online convex games. In *Proceedings of the nineteenth annual conference on computational learning theory*, 2008.
- [4] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. In *Conference on Learning Theory (COLT)*, 2010.
- [5] H. Brendan McMahan and Matthew Streeter. Adaptive bound optimization for online convex optimization. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, 2010.
- [6] Francesco Orabona and Dávid Pál. Coin betting and parameter-free online learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 577–585. Curran Associates, Inc., 2016.
- [7] Ashok Cutkosky and Kwabena Boahen. Online learning without prior information. *arXiv preprint arXiv:1703.02629*, 2017.
- [8] Tim van Erven and Wouter M Koolen. Metagrad: Multiple learning rates in online learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3666–3674. Curran Associates, Inc., 2016.
- [9] Wouter M Koolen, Peter Grünwald, and Tim van Erven. Combining adversarial guarantees and stochastic fast rates in online learning. In *Advances in Neural Information Processing Systems*, pages 4457–4465, 2016.
- [10] Francesco Orabona. Dimension-free exponentiated gradient. In *Advances in Neural Information Processing Systems*, pages 1806–1814, 2013.
- [11] Ashok Cutkosky and Kwabena A Boahen. Online convex optimization with unconstrained domains and losses. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 748–756. Curran Associates, Inc., 2016.
- [12] Brendan McMahan and Matthew Streeter. No-regret algorithms for unconstrained online convex optimization. In *Advances in neural information processing systems*, pages 2402–2410, 2012.
- [13] Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2):169–192, 2007.
- [14] Peter L Bartlett, Elad Hazan, and Alexander Rakhlin. Adaptive online gradient descent. In *NIPS*, volume 20, pages 65–72, 2007.
- [15] Amir Sani, Gergely Neu, and Alessandro Lazaric. Exploiting easy data in online optimization. In *Advances in Neural Information Processing Systems*, pages 810–818, 2014.

## A Theorems from Literature

In this section we reproduce here some previous theorems and notation for reference.

### A.1 Follow-the-Regularized-Leader

The Follow-the-Regularized-Leader (FTRL) framework for online optimization suggests choosing  $w_{t+1}$  according to the rule:

$$w_{t+1} = \underset{W}{\operatorname{argmin}} g_{1:t} \cdot w + \psi_t(w)$$

where  $\psi_t(w)$  is a function chosen by the algorithm called a *regularizer*. We use the following bound on the regret of FTRL, which is proved in [7]:

**Theorem 12.** *Let  $g_t, \dots, g_T$  be an arbitrary sequence of subgradients. Define  $g_0 = 0$  for notational convenience. Let  $\psi_0, \psi_1, \dots, \psi_{T-1}$  be a sequence of regularizer functions, such that  $\psi_t$  is chosen without knowledge of  $g_{t+1}, \dots, g_T$ . Let  $\psi_1^+, \dots, \psi_T^+$  be an arbitrary sequences of regularizer functions (possibly chosen with knowledge of the full subgradient sequence). Define  $w_1, \dots, w_T$  to be the outputs of FTRL with regularizers  $\psi_t$ :  $w_{t+1} = \operatorname{argmin} \psi_t(w) + g_{1:t} \cdot w$ , and define  $w_t^+$  for  $t = 2, \dots, T+1$  by  $w_{t+1}^+ = \operatorname{argmin} \psi_t^+(w) + g_{1:t} \cdot w$ . Then FTRL with regularizers  $\psi_t$  obtains regret*

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - u) &\leq \psi_T^+(u) - \psi_0(w_2^+) + \sum_{t=1}^T \psi_{t-1}(w_{t+1}^+) - \psi_t^+(w_{t+1}^+) + g_t \cdot (w_t - w_{t+1}^+) \\ &\quad + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) \end{aligned}$$

In the next subsection we recall the notion of an *adaptive regularizer* [7], which is a function  $\psi$  whose properties make it an easy building block for FTRL regularizers  $\psi_t$ . The analysis of FREEREXMOMENTUM is based upon the observation that its regularizers are constructed using an adaptive regularizer.

### A.2 Adaptive Regularizers

Before defining adaptive regularizers, we briefly introduce a minor generalization of strong-convexity below:

**Definition 13.** *Let  $W$  be a convex space and let  $\sigma : W^2 \rightarrow \mathbb{R}$  by an arbitrary function. We say a convex function  $f : W \rightarrow \mathbb{R}$  is  $\sigma(\cdot, \cdot)$ -strongly convex with respect to a norm  $\|\cdot\|$  if for all  $x, y \in W$  and  $g \in \partial f(x)$  we have*

$$f(y) \geq f(x) + g \cdot (y - x) + \frac{\min(\sigma(x), \sigma(y))}{2} \|x - y\|^2$$

We will exclusively make use of the special case  $\sigma(w, z) = \min(\sigma(w), \sigma(z))$ , and we will write  $\sigma$ -strongly convex instead of  $\sigma(\cdot)$ -strongly convex in all cases. Next we give the definition of adaptive regularizers:

**Definition 14.** *Any differentiable function  $\psi : W \rightarrow \mathbb{R}$  is called a  $(\sigma, \|\cdot\|)$ -adaptive regularizer if it that satisfies the following conditions:*

1.  $\psi(0) = 0$ .
2.  $\psi(x)$  is  $\sigma$ -strongly-convex with respect to some norm  $\|\cdot\|$  for some  $\sigma : W \rightarrow \mathbb{R}$  such that  $\|x\| \geq \|y\|$  implies  $\sigma(x) \leq \sigma(y)$ .
3. For any  $C$ , there exists a  $B$  such that  $\psi(x)\sigma(x) \geq C$  for all  $\|x\| \geq B$ .

Associated to every adaptive regularizer  $\psi$ , we define the function  $h(w) = \psi(w)\sigma(w)$ , and define  $h^{-1}(x) = \max_{h(x) \leq x} \|x\|$

Finally, we provide a general construction that converts an adaptive regularizer into a sequence of regularizers  $\psi_t$  used in FTRL (and in particular in FREEREXMOMENTUM). In the following we make use of the *dual norm*  $\|\cdot\|_*$ , which is defined by  $\|x\|_* = \sup_{\|y\|=1} x \cdot y$ .

**Definition 15.** Let  $\|\cdot\|$  be a norm and  $\|\cdot\|_*$  be the dual norm ( $\|x\|_* = \sup_{\|y\|=1} x \cdot y$ ). Let  $g_1, \dots, g_T$  be a sequence of subgradients and set  $L_t = \max_{t' \leq t} \|g_{t'}\|_*$ . Define the sequences  $\frac{1}{\eta_t}$  and  $a_t$  recursively by:

$$\begin{aligned} \frac{1}{\eta_0^2} &= 0 \\ \frac{1}{\eta_t^2} &= \max \left( \frac{1}{\eta_{t-1}^2} + 2\|g_t\|_*^2, L_t \|g_{1:t}\|_* \right) \\ a_1 &= \frac{1}{(L_1 \eta_1)^2} \\ a_t &= \max \left( a_{t-1}, \frac{1}{(L_t \eta_t)^2} \right) \end{aligned}$$

Suppose  $\psi$  is a  $(\sigma, \|\cdot\|)$ -adaptive regularizer and  $k > 0$ . Let  $\bar{w}_1, \dots, \bar{w}_T$  be an arbitrary sequence of vectors. Define

$$\begin{aligned} \psi_t(w) &= \frac{k}{\eta_t a_t} \psi(a_t(w - \bar{w}_t)) \\ w_{t+1} &= \operatorname{argmin}_{w \in W} \psi_t(w) + g_{1:t} \cdot w \end{aligned}$$

In order to use Theorem 12, we'll need to define some "shadow regularizers"  $\psi_t^+$ , which we do below:

**Definition 16.** Given a norm  $\|\cdot\|$  and a sequence of subgradients  $g_1, \dots, g_T$ , define  $L_t$  and  $\frac{1}{\eta_t}$  as in Definition 15, and define  $L_0 = L_1$ . We define  $\frac{1}{\eta_t^+}$  recursively by:

$$\begin{aligned} \frac{1}{\eta_0^+} &= \frac{1}{\eta_0} \\ \frac{1}{(\eta_t^+)^2} &= \max \left( \frac{1}{\eta_{t-1}^2} + 2\|g_t\|_* \min(\|g_t\|_*, L_{t-1}), L_{t-1} \|g_{1:t}\|_* \right) \end{aligned}$$

Further, given a  $k \geq 1$  and a non-decreasing sequence of positive numbers  $a_t$ , define  $\psi_t^+$  by:

$$\begin{aligned} \psi_t^+(w) &= \frac{k}{\eta_t^+ a_{t-1}} \psi(a_{t-1}(w - \bar{w}_{t-1})) \\ w_{t+1}^+ &= \operatorname{argmin}_{w \in W} \psi_t^+(w) + g_{1:t} \cdot w \end{aligned}$$

The following is the key technical Lemma from [7]. That paper does not take into account the "shifting" parameter  $\bar{w}_t$  and so technically the Lemma as proven there does not apply. However, by applying the change-of-coordinates  $w \mapsto w - \bar{w}_{t-1}$  we see that the "shifting" does not effect the conclusion.

**Lemma 17.** Suppose  $\psi$  is a  $(\sigma, \|\cdot\|)$ -adaptive regularizer and  $g_1, \dots, g_T$  is an arbitrary sequence of subgradients (possibly chosen adaptively). We use the regularizers of Definition 15. Recall that we define  $h(w) = \psi(w)\sigma(w)$  and  $h^{-1}(x) = \operatorname{argmax}_{h(w) \leq x} \|w\|$ . Define

$$\sigma_{\min} = \inf_{\|w\| \leq h^{-1}(10/k^2)} k\sigma(w)$$

and

$$D = 2 \max_t \frac{h^{-1} \left( 5 \frac{L_t}{k L_{t-1}} \right)}{a_{t-1}}$$

Then

$$\begin{aligned} &\psi_{t-1}(w_{t+1}^+) - \psi_t^+(w_{t+1}^+) + g_t(w_t - w_{t+1}^+) \\ &\leq \begin{cases} \|g_t\|_* \min(D, \max_t(\|w_t - w_{t+1}^+\|)) & \text{when } \|g_t\| > 2L_{t-1} \\ \frac{3\|g_t\|_*^2 \eta_t^+}{a_{t-1} \sigma_{\min}} & \text{otherwise} \end{cases} \end{aligned}$$

We copy over four final Lemmas from [7] that we include here for reference:

**Proposition 18.** Suppose  $\psi : W \rightarrow \mathbb{R}$  is a  $(\sigma, \|\cdot\|)$ -adaptive regularizer. Then  $\frac{\psi(aw)}{a}$  is an increasing function of  $a$  for all  $a > 0$  for all  $w \in W$ .

**Lemma 19.** Let  $\alpha_t$  be defined by

$$\alpha_0 = \frac{1}{(L_1 \eta_1)^2}$$

$$\alpha_t = \max \left( \alpha_{t-1}, \frac{1}{(L_t \eta_t)^2} \right)$$

Then

$$\frac{2(\|g\|_{\star})_{1:t}}{L_t} \geq a_t \geq \frac{2(\|g\|_{\star}^2)_{1:t}}{L_t^2}$$

**Lemma 20.** 1.

$$\sum_{t \mid \|g_t\|_{\star} \leq 2L_{t-1}} \|g_t\|_{\star}^2 \eta_t^+ \leq \frac{2}{\eta_T^+}$$

2. Suppose  $\alpha_t$  is defined by

$$\alpha_0 = \frac{1}{(L_1 \eta_1)^2}$$

$$\alpha_t = \max \left( \alpha_{t-1}, \frac{1}{(L_t \eta_t)^2} \right)$$

then

$$\sum_{t \mid \|g_t\|_{\star} \leq 2L_{t-1}} \|g_t\|_{\star}^2 \frac{\eta_t^+}{\alpha_{t-1}} \leq 15L_{\max}$$

**Lemma 21.** Let  $a_1, \dots, a_M$  be a sequence of non-negative numbers such that  $a_{i+1} \geq 2a_i$ . Then

$$\sum_{i=1}^M a_i \leq 2a_M$$

## B Proof of Main Theorem

### B.1 Proposition 5

First, we prove the simple Proposition 5, restated below for reference:

**Proposition 5.** If  $a, b, c$  and  $d$  are non-negative constants such that

$$x \leq a\sqrt{bx+c} + d$$

Then

$$x \leq 4a^2b + 2a\sqrt{c} + 2d$$

*Proof.* Suppose  $x \geq 2d$ . Then we have

$$\frac{x}{2} \leq a\sqrt{bx+c}$$

$$x^2 \leq 4a^2bx + 4a^2c$$

Now we use the quadratic formula to obtain

$$x \leq \frac{4a^2b}{2} + \frac{\sqrt{16a^4b^2 + 16a^2c}}{2}$$

$$\leq 4a^2b + 2a\sqrt{c}$$

Since we assumed  $x \geq 2d$  to obtain this bound, we conclude that  $x$  is at most the maximum of  $4a^2b + 2a\sqrt{c}$  and  $2d$ , which is bounded by their sum.  $\square$

### B.2 Proof of Theorem 1

Our strategy is based on the observation that FREEREXMOMENTUM is FTRL with regularizers  $\psi_t(w) = \frac{k}{a_t \eta_t} \phi(a_t \|w - \bar{w}_t\|)$  for  $\phi(x) = (x+1) \log(x+1) - x$  and  $k = \sqrt{5}$ , as can be easily verified by inspection of the updates. We will derive results for the case of arbitrary  $k$  and  $\bar{w}_t = \frac{\sum_{t'=1}^t \|g_{t'}\|_{w_{t'}}}{\delta + \|g\|_{1:t}}$  for arbitrary  $\delta$ , and then substitute  $k = \sqrt{5}$  and  $\delta = 1$  at the end to derive the bound for FREEREXMOMENTUM. We think this strategy clarifies the roles of the constants in the regret bound.

The following Theorem is nearly identical to the result in [7], but is very slightly generalized to our purposes:

**Theorem 22.** Suppose  $\psi$  is a  $(\sigma, \|\cdot\|)$ -adaptive regularizer and  $g_1, \dots, g_T$  is some arbitrary sequence of subgradients.

Set

$$\begin{aligned}\sigma_{\min} &= \inf_{\|w\| \leq h^{-1}(10/k^2)} k\sigma(w) \\ D &= \max_t \max \left( \frac{L_{t-1}^2}{(\|g\|_*^2)_{1:t-1}} h^{-1} \left( \frac{5L_t}{k^2 L_{t-1}} \right), \|w_t - w_{t+1}^+\| \right) \\ Q_T &= 2 \frac{\|g\|_{1:T}}{L_{\max}}\end{aligned}$$

Then FTRL with regularizers  $\psi_t$  achieves regret

$$R_T(u) \leq \frac{k}{Q_T \eta_T} \psi(Q_T(u - \overline{w_T})) + \frac{45L_{\max}}{\sigma_{\min}} + 2L_{\max}D + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+)$$

*Proof.* Using Theorem 12 and Lemma 17, our regret is bounded by

$$\begin{aligned}R_T(u) &\leq \psi_T(u) + \sum_{t=1}^T \psi_{t-1}(w_{t+1}^+) - \psi_t^+(w_{t+1}^+) + g_t(w_t - w_{t+1}^+) \\ &\quad + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) \\ &\leq \psi_T(u) + \sum_{t=1}^T \psi_{t-1}(w_{t+1}^+) - \psi_t^+(w_{t+1}^+) + g_t(w_t - w_{t+1}^+) \\ &\quad + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^T \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) \\ &\leq \psi_T(u) + \sum_{\|g_t\|_* \leq 2L_{t-1}} \frac{3\|g_t\|^2 \eta_t^+}{a_{t-1} \sigma_{\min}} + \sum_{\|g_t\|_* > 2L_{t-1}} \|g_t\|_* D' \\ &\quad + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+)\end{aligned}$$

where  $D'$  is defined by

$$D' = 2 \max_t \frac{h^{-1} \left( \frac{5L_t}{kL_{t-1}} \right)}{a_{t-1}}$$

Now we use Lemma 19 to conclude that

$$D' \leq D = \max_t \frac{L_{t-1}^2}{(\|g\|_*^2)_{1:t-1}} h^{-1} \left( \frac{5L_t}{kL_{t-1}} \right)$$

so that we have

$$\begin{aligned}R_T(u) &\leq \psi_T(u) + \sum_{\|g_t\|_* \leq 2L_{t-1}} \frac{3\|g_t\|^2 \eta_t^+}{a_{t-1} \sigma_{\min}} + \sum_{\|g_t\|_* > 2L_{t-1}} \|g_t\|_* D \\ &\quad + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^T \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+)\end{aligned}$$

Now using Lemma 20 we can simplify this to

$$\begin{aligned}R_T(u) &\leq \frac{k}{a_T \eta_T^+} \psi(a_T u) + \frac{45L_{\max}}{\sigma_{\min}} + \sum_{\|g_t\|_* > 2L_{t-1}} \|g_t\|_* D \\ &\quad + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^T \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+)\end{aligned}$$

Next, observe that each value of  $\|g_t\|_*$  in the sum  $\sum_{\|g_t\|_* > 2L_{t-1}} \|g_t\|_* D$  is at least twice the previous value, so that by Lemma 21 we conclude

$$\begin{aligned} R_T(u) &\leq \frac{k}{a_T \eta_T^+} \psi(a_T u) + \frac{45L_{\max}}{\sigma_{\min}} + 2L_{\max} D \\ &\quad + \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) \end{aligned}$$

Finally, we observe that (by Lemma 19),  $a_T \leq 2 \frac{\|g\|_{1:T}}{L_T} = Q_T$ , which gives the first inequality in the Theorem statement.  $\square$

We need the next theorem to convert  $\frac{45L_{\max}}{\sigma_{\min}}$  to  $405L_{\max}$ :

**Lemma 23.** *Suppose  $\psi(w) = ((\|w\| + 1) \log(\|w\| + 1) - \|w\|)$ . Then  $\psi$  is a  $(\frac{1}{\|\cdot\|+1}, \|\cdot\|)$ -adaptive regularizer. Using the terminology of Theorem 22, for  $k = \sqrt{5}$ ,  $\frac{45L_{\max}}{\sigma_{\min}} \leq 405L_{\max}$ .*

*Proof.* The fact that  $\psi$  is an adaptive regularizer is proved in [7] Proposition 9. For the second statement, we have

$$\begin{aligned} \frac{45L_{\max}}{\sigma_{\min}} &= \frac{45L_{\max}}{\inf_{\|w\| \leq h^{-1}(10/k^2)} k\sigma(w)} \\ &= \sup_{\|w\| \leq h^{-1}(10/k^2)} \frac{45L_{\max}(\|w\| + 1)}{k} \\ &= \frac{45L_{\max}(h^{-1}(10/k^2) + 1)}{k} \end{aligned}$$

Now it remains to compute an expression for  $h^{-1}$ . First we compute a bound on  $h$ :

$$\begin{aligned} h(w) &= \left( \log(\|w\| + 1) - \frac{\|w\|}{\|w\| + 1} \right) \\ &\geq \log(\|w\| + 1) - 1 \end{aligned}$$

so that

$$h^{-1}(x/k^2) \leq \exp(x/k^2 + 1) - 1$$

Now we numerically evaluate  $\frac{45L_{\max}}{\sigma_{\min}} = \frac{45L_{\max}(h^{-1}(10/k^2)+1)}{k}$  using  $k = \sqrt{5}$  to conclude the desired bound.  $\square$

So now we go to work to bound  $\psi_T^+(u) - \psi_T(u) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+)$ .

**Lemma 24.** *For any increasing sequence of numbers  $\{x_t\}$ ,*

$$\sum_{t=1}^T \frac{x_t - x_{t-1}}{x_t} \leq \log\left(\frac{x_T}{x_1}\right)$$

*Proof.* By concavity of log, we have

$$\log(x_t) - \log(x_{t-1}) \geq \frac{x_t - x_{t-1}}{x_t}$$

from which the result easily follows by telescoping a sum.  $\square$

**Lemma 25.** *Suppose  $\{x_t\}$  and  $\{\sigma_t\}$  are non-negative real numbers such that  $\sqrt{x_t}\sigma_t \geq \sqrt{x_{t-1}}\sigma_{t-1}$  for all  $t$ . Then*

$$\sum_{t=1}^T \frac{(x_t - x_{t-1})\sigma_t}{\sqrt{x_t}} \leq \sqrt{x_T}\sigma_T \log\left(\frac{x_T}{x_1}\right)$$

*Proof.* We have  $\sqrt{x_t}\sigma_t \leq \sqrt{x_T}\sigma_T$  so that

$$\sigma_t \leq \frac{\sqrt{x_T}\sigma_T}{\sqrt{x_t}}$$

Therefore

$$\begin{aligned} \sum_{t=1}^T \frac{(x_t - x_{t-1})\sigma_t}{\sqrt{x_t}} &\leq \sum_{t=1}^T \frac{(x_t - x_{t-1})\sqrt{x_T}\sigma_T}{x_t} \\ &\leq \sqrt{x_T}\sigma_T \log\left(\frac{x_T}{x_1}\right) \end{aligned}$$

□

We make a suggestive definition:

**Definition 26.** Given some  $\delta > 0$ ,

$$\begin{aligned} x_t &= \delta + (\|g\|_\star)_{1:t} \\ \bar{w}_t &= \frac{(\|g\|_\star w)_{1:t}}{x_t} \\ \sigma_t &= \sqrt{\frac{\delta\|\bar{w}_t\|^2 + \sum_{t'=1}^t \|g_{t'}\|_\star \|w_{t'} - \bar{w}_t\|^2}{x_t}} \end{aligned}$$

Observe that the values of  $\bar{w}_t$  given in the psuedo-code for FREEREXMOMENTUM match the values above for  $\delta = 1$ . We will carry through all our calculations for general  $\delta$ , and then substitute  $\delta = 1$  at the very end to obtain our regret bound.

Consider a random vector that takes on value  $w_t \neq 0$  for  $t \leq T$  with probability proportional to  $\|g_t\|_\star$  and value 0 with probability proportional to  $\delta + \sum_{w_t=0}^T \|g_t\|_\star$ . Then the expectation of this vector is  $\bar{w}_T$  and  $\sigma_T^2$  is its variance. Thus for any vector  $X$ , by a standard bias-variance decomposition we have

$$\delta X^2 + \sum_{t=1}^T \|g_t\|_\star \|X - \bar{w}_T\|^2 = x_T(\sigma_T^2 + \|X - \bar{w}_T\|^2)$$

**Lemma 27.** Using the definitions in Definition 26, for all  $T$ :

$$\sigma_T \sqrt{x_T} - \sigma_{T-1} \sqrt{x_{T-1}} \geq \frac{\|g_T\|_\star \|w_T - \bar{w}_T\|^2}{2\sigma_T \sqrt{x_T}}$$

*Proof.*

$$\begin{aligned} \sigma_T \sqrt{x_T} &= \sqrt{\delta\|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\|_\star \|w_t - \bar{w}_T\|^2} \\ &\geq \sqrt{\delta\|\bar{w}_T\|^2 + \sum_{t=1}^{T-1} \|g_t\|_\star \|w_t - \bar{w}_T\|^2} + \frac{\|g_T\|_\star \|w_T - \bar{w}_T\|^2}{2\sqrt{\delta\|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\|_\star \|w_t - \bar{w}_T\|^2}} \\ &= \sqrt{\delta\|\bar{w}_T\|^2 + \sum_{t=1}^{T-1} \|g_t\|_\star \|w_t - \bar{w}_T\|^2} + \frac{\|g_T\|_\star \|w_T - \bar{w}_T\|^2}{2\sigma_T \sqrt{x_T}} \end{aligned}$$

And also we have

$$\begin{aligned} \delta\|\bar{w}_T\|^2 + \sum_{t=1}^{T-1} \|g_t\|_\star \|w_t - \bar{w}_T\|^2 &= x_{T-1}(\sigma_{T-1}^2 + \|\bar{w}_T - \bar{w}_{T-1}\|^2) \\ &\geq x_{T-1}\sigma_{T-1}^2 \end{aligned}$$

and so we can conclude the desired inequality. □

**Lemma 28.** Again using the terms from Definition 26, we have

$$\sum_{t=1}^T \frac{\|g_t\|_\star \|w_t - \bar{w}_t\|}{\sqrt{x_t}} \leq \sigma_T \sqrt{x_T} \left( 2 + \log\left(\frac{x_T}{\delta + L_1}\right) \right)$$

*Proof.* From Lemma 27, we see that when  $\|w_t - \bar{w}_t\| \geq \sigma_t$ , we have  $\frac{\|g_t\|_* \|w_t - \bar{w}_t\|}{\sqrt{x_t}} \leq 2\sigma_t \sqrt{x_t} - 2\sigma_{t-1} \sqrt{x_{t-1}}$  so that we can write:

$$\sum_{t=1}^T \frac{\|g_t\|_* \|w_t - \bar{w}_t\|}{\sqrt{x_t}} \leq 2\sigma_T \sqrt{x_T} + \sum_{t=1}^T \frac{\|g_t\|_* \sigma_t}{\sqrt{x_t}}$$

Now we observe (e.g. by Lemma 27) that  $\sigma_t \sqrt{x_t} \geq \sigma_{t-1} \sqrt{x_{t-1}}$  for all  $t$  and that  $\|g_t\|_* = x_t - x_{t-1}$  so that applying Lemma 25 we have

$$\sum_{t=1}^T \frac{\|g_t\|_* \|w_t - \bar{w}_t\|}{\sqrt{x_t}} \leq 2\sigma_T \sqrt{x_T} + \sigma_T \sqrt{x_T} \log \left( \frac{x_T}{x_1} \right)$$

as desired.  $\square$

**Proposition 29.** *Let  $a_1, \dots, a_T$  be non-negative numbers. Then*

$$\sum_{t=1}^T \frac{a_t}{(a_{1:t})^{3/2}} \leq \frac{3}{\sqrt{a_1}} - \frac{2}{\sqrt{a_{1:T}}}$$

*Proof.* We proceed by induction. For the base case, we have

$$\sum_{t=1}^1 \frac{a_t}{(a_{1:t})^{3/2}} = \frac{1}{\sqrt{a_1}}$$

Suppose that  $\sum_{t=1}^T \frac{a_t^2}{(a_{1:t})^{3/2}} \leq \frac{3}{\sqrt{a_1}} - \frac{2}{\sqrt{a_{1:T}}}$ .

By concavity of  $-\frac{1}{\sqrt{x}}$  we have

$$\left( \frac{3}{\sqrt{a_1}} - \frac{2}{\sqrt{a_{1:T+1}}} \right) - \left( \frac{3}{\sqrt{a_1}} - \frac{2}{\sqrt{a_{1:T}}} \right) \geq \frac{a_{T+1}}{(a_{1:T+1})^{3/2}}$$

By the induction assumption we have

$$\begin{aligned} \sum_{t=1}^{T+1} \frac{a_t}{(a_{1:t})^{3/2}} &\leq \frac{3}{\sqrt{a_1}} - \frac{2}{\sqrt{a_{1:T}}} + \frac{a_{T+1}}{(a_{1:T+1})^{3/2}} \\ &\leq \frac{3}{\sqrt{a_1}} - \frac{2}{\sqrt{a_{1:T+1}}} \end{aligned}$$

as desired.  $\square$

**Lemma 30.** *Define  $\bar{w}_t$  as in Definition 26. Define  $M_t = \sup_{w, w' \in W} \|\nabla \psi(a_t(w - w'))\|_*$ . Then using the terminology of Definition 15, we have*

$$\begin{aligned} \psi_T^+(u) - \psi_T(u) + \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) &\leq \sigma_T \sqrt{2L_{\max} x_T} \left( 2 + \log \left( \frac{x_T}{x_1} \right) \right) \max_t M_t \\ &\quad + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{\delta} + L_1} \max_t \|\bar{w}_{t-1} - w_t\| \max_t M_t \end{aligned}$$

*Proof.* From Proposition 18, we see that  $\frac{1}{a_{t-1}} \psi(a_{t-1} x) \leq \frac{1}{a_t} \psi(a_t x)$  for all  $x$ . Therefore we have:

$$\begin{aligned} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) &= \frac{1}{\eta_t^+ a_{t-1}} \psi(a_{t-1}(w_{t+2}^+ - \bar{w}_{t-1})) - \frac{1}{\eta_t a_t} \psi(a_t(w_{t+2}^+ - \bar{w}_t)) \\ &\leq \frac{1}{\eta_t a_t} \psi(a_t(w_{t+2}^+ - \bar{w}_{t-1})) - \frac{1}{\eta_t a_t} \psi(a_t(w_{t+2}^+ - \bar{w}_t)) \\ &\leq \frac{1}{\eta_t a_t} \|\nabla \psi(a_t(w_{t+2}^+ - \bar{w}_{t-1}))\|_* a_t \|\bar{w}_t - \bar{w}_{t-1}\| \\ &\leq \frac{\|\bar{w}_t - \bar{w}_{t-1}\|}{\eta_t} \max_t \|\nabla \psi(a_t(w_{t+2}^+ - \bar{w}_{t-1}))\|_* \\ &\leq \|\bar{w}_t - \bar{w}_{t-1}\| \sqrt{2L_{\max} x_t} \max_t \|\nabla \psi(a_t(w_{t+2}^+ - \bar{w}_{t-1}))\|_* \\ &\leq \|\bar{w}_t - \bar{w}_{t-1}\| \sqrt{2L_{\max} x_t} \max_t M_t \end{aligned}$$



Where in the last step we observe  $\frac{1}{\eta_t} \leq \sqrt{2L_{\max}(\|g\|_*)_{1:t}} \leq \sqrt{2L_{\max}x_t}$ , which can be easily deduced by induction, or from Proposition 19 of [7].

The exact same argument can be used to show

$$\psi_T^+(u) - \psi_T(u) \leq \|\bar{w}_T - \bar{w}_{T-1}\| \sqrt{2L_{\max}x_T} \max_t M_t$$

Next we characterize  $\bar{w}_t - \bar{w}_{t-1}$ :

$$\begin{aligned} \bar{w}_{t-1} - \bar{w}_t &= \bar{w}_{t-1} - \frac{(\delta + (\|g\|_*)_{1:t-1})\bar{w}_{t-1} + \|g_t\|_* w_t}{\delta + (\|g\|_*)_{1:t}} \\ &= \frac{\|g_t\|_*}{\delta + (\|g\|_*)_{1:t}} (\bar{w}_{t-1} - w_t) \end{aligned}$$

We can take this calculation one step further:

$$\begin{aligned} \psi_T^+(u) - \psi_T(u) + \bar{w}_{t-1} - \bar{w}_t &= \frac{\|g_t\|_*}{\delta + (\|g\|_*)_{1:t}} (\bar{w}_{t-1} - w_t) \\ &= \frac{\|g_t\|_*}{\delta + (\|g\|_*)_{1:t}} (\bar{w}_t - w_t) + \frac{\|g_t\|_*}{\delta + (\|g\|_*)_{1:t}} (\bar{w}_{t-1} - \bar{w}_t) \\ &= \frac{\|g_t\|_*}{\delta + (\|g\|_*)_{1:t}} (\bar{w}_t - w_t) + \frac{\|g_t\|_*^2}{(\delta + (\|g\|_*)_{1:t})^2} (\bar{w}_{t-1} - w_t) \end{aligned}$$

Thus we have

$$\begin{aligned} \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) &\leq \sum_{t=1}^T \|\bar{w}_t - \bar{w}_{t-1}\| \sqrt{2L_{\max}x_t} \max_t M_t \\ &= \sum_{t=1}^T \frac{\sqrt{2L_{\max}} \|g_t\|_* \|\bar{w}_t - w_t\|}{\sqrt{x_t}} \max_t M_t \\ &\quad \sum_{t=1}^T \frac{\sqrt{2L_{\max}} \|g_t\|_*^2 \|\bar{w}_{t-1} - w_t\|}{x_t^{3/2}} \max_t M_t \\ &\leq \sigma_T \sqrt{2L_{\max}x_T} \left( 2 + \log \left( \frac{x_T}{x_1} \right) \right) \max_t M_t \\ &\quad + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{\delta + L_1}} \max_t \|\bar{w}_{t-1} - w_t\| \max_t M_t \end{aligned}$$

Where we've used Proposition 29 to conclude that

$$\sum_{t=1}^T \frac{\|g_t\|_*^2}{x_t^{3/2}} \leq \frac{3L_{\max}}{\sqrt{\delta + L_1}}$$

and also used Lemma 28 in the last inequality.  $\square$

Now if we restrict ourselves to a bounded domain of diameter  $B$  and use the regularizer  $\psi(w) = (\|w\| + 1) \log(\|w\| + 1) - \|w\|$ , we obtain

$$\max_t M_t \leq \log(Ba_T)$$

so that we have

$$\begin{aligned} \sum_{t=1}^{T-1} \psi_t^+(w_{t+2}^+) - \psi_t(w_{t+2}^+) &\leq \sqrt{2L_{\max} \left( \delta \|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\| \|w_t - \bar{w}_T\|^2 \right)} \left( 2 + \log \left( \frac{\delta + \|g\|_{1:T}}{\delta + \|g_1\|} \right) \right) \log(Ba_T) \\ &\quad + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{\delta + L_1}} B \log(Ba_T) \end{aligned}$$

Combining this with Theorem 22 and Lemma 23 and using  $\delta = 1$  and  $k = \sqrt{5}$  we have proved a regret bound on FTRL with regularizers  $\psi_t = \frac{\sqrt{5}}{\eta_t} \psi(w_t - \bar{w}_t)$  with  $\psi = (\|w\| + 1) \log(\|w\| + 1) - \|w\|$ . Recall that FREEREXMOMENTUM is precisely FTRL with these regularizers, so we have proved Theorem 1:

**Theorem 1.** Let  $\psi(w) = (\|w\| + 1) \log(\|w\| + 1) - \|w\|$ . Set  $L_t = \max_{t' \leq t} \|g_{t'}\|$ , and  $Q_T = 2 \frac{\|g\|_{1:T}}{L_{\max}}$ . Define  $\frac{1}{\eta_t}$  and  $\alpha_t$  as in the pseudo-code for FREEREXMOMENTUM (Algorithm 1). Then the regret of FREEREXMOMENTUM is bounded by:

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq \frac{\sqrt{5}}{Q_T \eta_T} \psi(Q_T(w^* - \bar{w}_T)) + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(Ba_T + 1) \\ &\quad + \sqrt{2L_{\max} \left( \|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\| \|w_t - \bar{w}_T\|^2 \right)} \left( 2 + \log \left( \frac{1 + \|g\|_{1:T}}{1 + \|g_1\|} \right) \right) \log(Ba_T + 1) \end{aligned}$$

### B.3 Proof of Corollaries 2 and 3

First we prove Corollary 2, restated below:

**Corollary 2.** Under the assumptions and notation of Theorem 1, the regret of FREEREXMOMENTUM is bounded by:

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right)} \log(2BT + 1)(2 + \log(T)) \\ &\quad + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(2BT + 1) \end{aligned}$$

*Proof.* We need the observations

$$\begin{aligned} \psi(w) &\leq \|w\| \log(\|w\| + 1) \\ \frac{1}{\eta_T} &\leq \sqrt{2L_{\max}(1 + \|g\|_{1:T})} \\ \alpha_T &\leq 2T \end{aligned}$$

Using these identities with Theorem 1 gives us

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq \sqrt{5} \sqrt{2\|w^* - \bar{w}_T\|^2 L_{\max}(1 + \|g\|_{1:T})} \log(2BT + 1) \\ &\quad + \sqrt{2L_{\max} \left( \|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\| \|w_t - \bar{w}_T\|^2 \right)} (2 + \log(T)) \log(2BT + 1) \\ &\quad + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(2BT + 1) \end{aligned}$$

Now use  $\sqrt{a} + \sqrt{b} \leq \sqrt{2a + 2b}$  to reach the conclusion:

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq 2\sqrt{5} \sqrt{L_{\max} \left( \|w^* - \bar{w}_T\|^2 (1 + \|g\|_{1:T}) + \|\bar{w}_T\|^2 + \sum_{t=1}^T \|g_t\| \|w_t - \bar{w}_T\|^2 \right)} \\ &\quad \times \log(2BT + 1)(2 + \log(T)) \\ &\quad + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(2BT + 1) \\ &\leq 2\sqrt{5} \sqrt{L_{\max} \left( \|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right)} \log(2TB + 1)(2 + \log(T)) \\ &\quad + 405L_{\max} + 2L_{\max}B + 3 \frac{L_{\max} \sqrt{2L_{\max}}}{\sqrt{1+L_1}} B \log(2BT + 1) \end{aligned}$$

□

Now we Corollary 3, again restated below:

**Corollary 3.** *The regret of coordinate-wise FREEREXMOMENTUM is bounded by:*

$$\begin{aligned} \sum_{t=1}^T g_t \cdot (w_t - w^*) &\leq 2\sqrt{5} \sqrt{dL_{\max} \left( d\|w^*\|^2 + \sum_{t=1}^T \|g_t\| \|w^* - w_t\|^2 \right)} \log(2Tb + 1)(2 + \log(T)) \\ &\quad + 405dL_{\max} + 2L_{\max}db + 3d \frac{L_{\max}\sqrt{2L_{\max}}}{\sqrt{1 + L_1}} b \log(2bT + 1) \end{aligned}$$

*Proof.* The Corollary follows by application of Cauchy-Schwarz inequality to Corollary 2. Recall that

$$R_T(u) \leq \sum_{t=1}^T g_t \cdot (w_t - u) = \sum_{i=1}^d \sum_{t=1}^T g_{t,i} (w_{t,i} - u_i)$$

So that the regret can be computed by summing the regret bound of Corollary 2 across dimensions:

$$\begin{aligned} R_T(u) &\leq 2\sqrt{5} \sum_{i=1}^d \sqrt{L_{\max} \left( (w_i^*)^2 + \sum_{t=1}^T |g_{t,i}| (w_i^* - w_{t,i})^2 \right)} \log(2bT + 1)(2 + \log(T)) \\ &\quad + d405L_{\max} + 2dL_{\max}b + 3d \frac{L_{\max}\sqrt{2L_{\max}}}{\sqrt{1 + L_1}} b \log(2bT + 1) \\ &\leq 2\sqrt{5} \sqrt{dL_{\max} \left( d\|w^*\|^2 + \sum_{i=1}^d \sum_{t=1}^T |g_{t,i}| (w_i^* - w_{t,i})^2 \right)} \log(2bT + 1)(2 + \log(T)) \\ &\quad + d405L_{\max} + 2dL_{\max}b + 3d \frac{L_{\max}\sqrt{2L_{\max}}}{\sqrt{1 + L_1}} b \log(2bT + 1) \\ &\leq 2\sqrt{5} \sqrt{dL_{\max} \left( d\|w^*\|^2 + \sum_{t=1}^T \|g_t\| \sqrt{\sum_{i=1}^d (w_i^* - w_{t,i})^4} \right)} \log(2bT + 1)(2 + \log(T)) \\ &\quad + d405L_{\max} + 2dL_{\max}b + 3d \frac{L_{\max}\sqrt{2L_{\max}}}{\sqrt{1 + L_1}} b \log(2bT + 1) \\ &\leq 2\sqrt{5} \sqrt{dL_{\max} \left( d\|w^*\|^2 + \sum_{t=1}^T \|g_t\| \sum_{i=1}^d \|w_i^* - w_{t,i}\|^2 \right)} \log(2bT + 1)(2 + \log(T)) \\ &\quad + d405L_{\max} + 2dL_{\max}b + 3d \frac{L_{\max}\sqrt{2L_{\max}}}{\sqrt{1 + L_1}} b \log(2bT + 1) \end{aligned}$$

where the first inequality follows from convexity of  $\sqrt{x}$ , the second from Cauchy-Schwarz, and the third because  $\|x\|_4^2 = \sqrt{\sum_{i=1}^d x_i^4} \leq \|x\|_2^2$ .  $\square$