

---

# No-Regret is not enough! Bandits with General Constraints through Adaptive Regret Minimization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In the bandits with knapsacks framework (BwK) the learner has  $m$  resource-  
2 consumption (*i.e.*, packing) constraints. We focus on the generalization of BwK in  
3 which the learner has a set of general long-term constraints. The goal of the learner  
4 is to maximize their cumulative reward, while at the same time achieving small  
5 cumulative constraints violations. In this scenario, there exist simple instances  
6 where conventional methods for BwK fail to yield sublinear violations of constraints.  
7 We show that it is possible to circumvent this issue by requiring the primal and dual  
8 algorithm to be *weakly adaptive*. Indeed, even in absence of any information on  
9 the Slater’s parameter  $\rho$  characterizing the problem, the interplay between weakly  
10 adaptive primal and dual regret minimizers yields a “self-bounding” property of  
11 dual variables. In particular, their norm remains suitably upper bounded across  
12 the entire time horizon even without explicit projection steps. By exploiting this  
13 property, we provide *best-of-both-worlds* guarantees for stochastic and adversarial  
14 inputs. In the first case, we show that the algorithm guarantees sublinear regret. In  
15 the latter case, we establish a tight competitive ratio of  $\rho/(1 + \rho)$ . In both settings,  
16 constraints violations are guaranteed to be sublinear in time. Finally, this results  
17 allow us to obtain new result for the problem of *contextual bandits with linear*  
18 *constraints*, providing the first no- $\alpha$ -regret guarantees for adversarial contexts.

## 19 1 Introduction

20 We consider a problem in which a decision maker tries to maximize their cumulative reward over  
21 a time horizon  $T$ , subject to a set of  $m$  *long-term constraints*. At each round  $t$ , the learner chooses  
22  $x_t \in \mathcal{X}$  and, subsequently, observes a reward  $f_t(x_t) \in [0, 1]$  and  $m$  constraint functions  $\mathbf{g}_t(x_t) \in$   
23  $[-1, 1]^m$ . Then, the problem becomes that of finding a sequence of decisions which guarantees a  
24 reward close to that of the best fixed decision in hindsight, while satisfying long-term constraints  
25  $\sum_{t=1}^T \mathbf{g}_t(\mathbf{x}_t) \leq \mathbf{0}$  up to small sublinear violations. This framework subsumes the *bandits with*  
26 *knapsacks* (BwK) problem, where there are only resource-consumption constraints [10, 5, 30].

27 Inputs  $(f_t, \mathbf{g}_t)$  may be either stochastic or adversarial. The goal is designing algorithms providing  
28 guarantees for both input models, without prior knowledge of the specific environment they will  
29 encounter. Achieving this goal involves addressing two crucial challenges which prevent a direct  
30 application of primal-dual approaches based on the LagrangeBwK framework in [30].

### 31 1.1 Technical Challenges

32 In order to obtain meaningful regret guarantees, primal-dual frameworks based on LagrangeBwK  
33 need to control the magnitude of dual variables. This is necessary as dual variables appear in the loss  
34 function of the primal algorithm, and, therefore, influence the no-regret guarantees provided by the

35 primal algorithm. In the context of knapsack constraints, this is usually achieved by exploiting the  
 36 existence of a strictly feasible solution with Slater’s parameter  $\rho$ , consisting of a *void action* which  
 37 yields zero reward and resource consumption. For instance, the frameworks of [14, 17] guarantee  
 38 boundedness of dual multipliers through an explicit projection step on the interval  $[0, 1/\rho]$ . However,  
 39 in settings with general constraints beyond resource consumption, it is often unreasonable to assume  
 40 that the learner knows the Slater’s parameter  $\rho$  a priori. The problem of operating without knowledge  
 41 of  $\rho$  has been already addressed in the stochastic setting [4, 5, 45, 44, 19]. For instance, a simple  
 42 approach for the case of stochastic inputs involves adding an initial estimation phase to calculate  
 43 an estimate of  $\rho$ , and subsequently treating this estimate as the true parameter [19]. However, these  
 44 techniques cannot be applied in adversarial environments as estimates of  $\rho$  based on the initial rounds  
 45 could be inaccurate about future inputs.

46 Primal-dual templates based on `LagrangeBwK` usually operate under the assumption that the  
 47 primal and dual algorithms have the no-regret property. In the case of standard `BwK`, the no-regret  
 48 requirement is sufficient to obtain optimal guarantees (see, e.g., [30, 17]). However, in our model,  
 49 there exist simple instances in which the primal and dual algorithms satisfy the no-regret requirement,  
 50 but the overall framework fails to guarantee small constraints violations (see Section 5.1). Moreover,  
 51 known techniques to prevent this problem, such as introducing a *recovery phase* to prevent excessive  
 52 violations, crucially require a priori knowledge of the Slater’s parameter  $\rho$  [19].

## 53 1.2 Contributions

54 Our approach is based on a generalization of the technique presented in [18] for online bidding under  
 55 one budget and one return-on-investments constraint. The crux of the approach is requiring that both  
 56 the primal and dual algorithms are *weakly adaptive*, that is, they guarantee a regret upper bound  
 57 of  $o(T)$  for each sub-interval of the time horizon [29]. We generalize this approach to the case of  
 58  $m$  general constraints, thereby providing the first primal-dual framework for this problem that can  
 59 operate without any knowledge of Slater’s parameter in both stochastic and adversarial environments.

60 First, we prove a “self-bounding” lemma for the case of  $m$  arbitrary constraints. It shows that, if the  
 61 primal and dual algorithms are weakly adaptive, then boundedness of dual multipliers emerges as a  
 62 byproduct of the interaction between the primal and dual algorithm. Thus, it is possible to guarantee  
 63 a suitable upper bound on the dual multipliers even without any information on Slater’s parameter.

64 We use this result to prove *best-of-both-worlds* no-regret guarantees for primal-dual frameworks  
 65 derived from `LagrangeBwK` which employ weakly adaptive primal and dual algorithms. Our  
 66 guarantees will be modular with respect to the regret guarantees of the primal and dual algorithms.  
 67 In presence of a suitable primal regret minimizer, we show that our framework yields the following  
 68 no-regret guarantees while attaining sublinear constraints violations: in the stochastic setting, it  
 69 guarantees sublinear regret with respect to the best fixed randomized strategy that is feasible in  
 70 expectation. Remarkably, this result is obtained without having to allocate the initial  $T^{1/2}$  rounds for  
 71 estimating the unknown parameter as in [19]. In the adversarial setting, our framework guarantees  
 72 a competitive ratio of  $\rho/(1 + \rho)$  against the best unconstrained strategy in hindsight. We provide a  
 73 lower bound showing that this cannot be improved if constraint violations have to be  $o(T)$ . This is  
 74 the first regret guarantee for our problem in adversarial environments.

75 Finally, we show that our model can be used to describe the *contextual bandits with linear constraints*  
 76 (`CBwLC`) problem, which was recently studied by [40, 27] in the context of stochastic and non-  
 77 stationary environments. Our framework allows to extend these works in two directions: we establish  
 78 the first no- $\alpha$ -regret guarantees for `CBwLC` when contexts are generated by an adversary, and we  
 79 provide the first  $\tilde{O}(\sqrt{T})$  guarantees for the stochastic setting when the learner does not know an  
 80 estimate of the Slater’s parameter of the problem.

## 81 2 Related Work

82 **Bandits with Knapsacks.** The (stochastic) `BwK` problem was introduced and optimally solved by  
 83 [9, 10]. Other algorithms with optimal regret guarantees have been proposed by [4, 5], whose  
 84 approach is based on the paradigm of *optimism in the face of uncertainty*, and in [31, 30]. In the  
 85 latter works, the authors propose the `LagrangeBwK` framework, which has a natural interpretation:  
 86 arms can be thought of as primal variables, and resources as dual variables. The framework works by

87 setting up a repeated two-player zero-sum game between a primal and a dual player, and by showing  
 88 convergence to a Nash equilibrium of the expected Lagrangian game.

89 **Adversarial BwK.** The adversarial BwK problem was first introduced in [31, 30], where they studied  
 90 the case in which the learner has  $m$  knapsack constraints, and inputs are selected by an oblivious  
 91 adversary. Their algorithm is based on a modified analysis of `LagrangeBwK`, and guarantees a  
 92  $O(m \log T)$  competitive ratio. Subsequently, [32] provided a new analysis obtaining a  $O(\log m \log T)$   
 93 competitive ratio, which is optimal. In the case in which budgets are  $\Omega(T)$ , [17] showed that it is  
 94 possible to achieve a constant competitive ratio of  $1/\rho$  where  $\rho$  is the per-iteration budget.

95 **Beyond packing constraints.** [17] studies a setting with general constraints analogous to ours, and  
 96 show how to adapt the `LagrangeBwK` framework to obtain best-of-both-worlds guarantees when  
 97 Slater’s parameter is known a priori. Similar guarantees are also provided, in the stochastic setting,  
 98 by [40], which then extend the results to the `CBwLC` model. Finally, the work of [18] introduces  
 99 the use of weakly adaptive regret minimizers within the `LagrangeBwK` framework, and provides  
 100 guarantees in the specific case of one budget constraint and one return-on-investments constraint.

101 **Contextual bandits (CB).** We briefly survey the most relevant works for our paper. Further references  
 102 can be found in [39, Chapter 8]. As in [41], we focus on CB with regression oracles [24, 25, 16, 38].  
 103 The contextual version of BwK was first studied by [11] in the case of classification oracles. A  
 104 regret-optimal and oracle-efficient algorithm for this problem was proposed by [6] by exploiting the  
 105 oracle-efficient algorithm for CB by [2]. The first regression-based approach for constrained BwK  
 106 was proposed by [3] by exploiting the optimistic approach for linear CB [34, 21, 1]. [27] propose a  
 107 regression-based approach for a constrained BwK setup under stochastic inputs. Finally, a notable  
 108 special case of constrained CB is online bidding under constraints [13, 20, 26, 22, 43].

109 **Other related works.** [23] show how to interpolate between the fully stochastic and the fully  
 110 adversarial setting, depending on the magnitude of fluctuations in expected rewards and consumptions  
 111 across rounds. [35] study a non-stationary setting and provide no-regret guarantees against the best  
 112 dynamic policy through a UCB-based algorithm. Some recent works explore the case in which  
 113 resource consumptions in BwK can be non-monotonic [33, 15]. Finally, a related line of works is the  
 114 one on online allocation problems with fixed per-iteration budget, where the input pair of reward and  
 115 costs is observed *before* the learner makes a decision [14, 12].

### 116 3 Preliminaries

117 There are  $T$  rounds and  $m$  constraints. We denote with  $\mathcal{X} \subset \mathbb{R}^K$  the decision space of the agent.  
 118 At each round  $t \in \llbracket T \rrbracket$ , the agent selects an action  $x_t \in \mathcal{X}$  and subsequently observes a reward  
 119  $f_t(x_t)$  and costs function  $g_t(x_t) \in [-1, 1]^m$ , with  $f_t : \mathcal{X} \rightarrow [0, 1]$  and  $g_{t,i} : \mathcal{X} \rightarrow [-1, 1]$  for  
 120 each  $i \in \llbracket m \rrbracket$ .<sup>1</sup> The reward and cost functions can either be chosen by an oblivious adversary or  
 121 drawn from a distribution. The goal of the decision maker is to maximize the cumulative reward  
 122  $\text{Rew}(T) := \sum_{t \in \llbracket T \rrbracket} f_t(x_t)$ , while minimizing the cumulative violation  $V_i(T)$  defined as

$$V_i(T) := \sum_{t \in \llbracket T \rrbracket} g_{t,i}(x_t)$$

123 for each constraint  $i \in \llbracket m \rrbracket$ . We denote by  $V(T) := \max_{i \in \llbracket m \rrbracket} V_i(T)$  the maximum cumulative  
 124 violation across the  $m$  constraints.

#### 125 3.1 Baselines

126 We will provide best-of-both-worlds no-regret guarantees for our algorithm, meaning that it achieves  
 127 optimal theoretical guarantees both in the stochastic and adversarial setting. In this section, we  
 128 introduce the baselines used to define the regret in these two scenarios.

129 **Adversarial Setting** In the adversarial setting we employ the strongest baseline possible, *i.e.*, the  
 130 best *unconstrained* strategy in hindsight:

$$\text{Opt}_{\text{Adv}} := \sup_{x \in \mathcal{X}} \sum_{t \in \llbracket T \rrbracket} f_t(x).$$

<sup>1</sup>In this work, for any  $a, b \in \mathbb{N}$ , with  $a < b$  we denote with  $\llbracket a \rrbracket$  the set  $\{1, \dots, a\}$  while  $\llbracket a, b \rrbracket$  the set  $\{a + 1, \dots, b\}$ .

131 This baseline is more powerful than the best fixed strategy which is feasible on average [30, 17],  
 132 which is the most common baseline in the literature. Our algorithm will yield an optimal competitive  
 133 ratio against this stronger baseline. In this setting, we define  $\rho_{\text{Adv}}$  as the feasibility parameter of the  
 134 problem instance, *i.e.*, the largest reduction of cumulative violations that the agent is guaranteed to  
 135 achieve by playing a “safe” strategy  $\xi^\circ \in \Delta(\mathcal{X})$ , where  $\Delta(\mathcal{X})$  is the set of all probability measures  
 136 on  $\mathcal{X}$ . Formally,

$$\rho_{\text{Adv}} := - \max_{t \in \llbracket T \rrbracket, i \in \llbracket m \rrbracket} \mathbb{E}_{x \sim \xi^\circ} [g_{t,i}(x)] \quad \text{and} \quad \xi^\circ := \arg \inf_{\xi \in \Delta(\mathcal{X})} \max_{t \in \llbracket T \rrbracket, i \in \llbracket m \rrbracket} \mathbb{E}_{x \sim \xi} [g_{t,i}(x)].$$

137 **Stochastic Setting** When the reward and the costs are stochastic we denote by  $\bar{f}$  and  $\bar{g}$  the mean of  
 138  $f_t$  and  $g_t$ , respectively. In particular, we have that the rewards are drawn so that  $\mathbb{E}_{\text{Env}}[f_t(x)] = \bar{f}(x)$   
 139 (and similarly for the costs), where  $\mathbb{E}_{\text{Env}}$  denotes expectation over the environment measure. We  
 140 define the baseline for the stochastic setting as the best fixed *randomized* strategy that satisfies the  
 141 constraints in expectation, which is the standard choice in Stochastic Bandits with Knapsacks settings  
 142 [9, 30]. Formally,

$$\text{Opt}_{\text{Stoc}} := \sup_{\xi \in \Delta(\mathcal{X}) : \mathbb{E}_{x \sim \xi} [\bar{g}(x)] \leq \mathbf{0}} \mathbb{E}_{x \sim \xi} [\bar{f}(x)].$$

143 Similarly to the adversarial case, we define the feasibility parameter  $\rho_{\text{Stoc}}$  as the “most negative”  
 144 cost achievable by randomized strategies *in expectation*:

$$\rho_{\text{Stoc}} := - \inf_{\xi \in \Delta(\mathcal{X})} \max_{i \in \llbracket m \rrbracket} \mathbb{E}_{x \sim \xi} [\bar{g}_i(x)].$$

145 As it is customary in relevant literature (see, *e.g.*, [30, 17, 19]), we make the following natural  
 146 assumption about the existence of a strictly feasible solution. Note that we do not make any  
 147 assumption on the variance of the samples ( $f_t, g_t$ ) as we assume that they have bounded support,  
 148 *i.e.*, with probability holds that  $f_t(x) \in [0, 1]$  and  $g_{t,i}(x) \in [-1, 1]$  for all  $x \in \mathcal{X}$  and  $i \in \llbracket m \rrbracket$ .

149 **Assumption 3.1.** In the adversarial setting, the sequence of inputs  $(f_t, g_t)_{t=1}^T$  is such that  $\rho_{\text{Adv}} > 0$ .  
 150 In the stochastic setting, the environment  $\text{Env}$  is such that  $\rho_{\text{Stoc}} > 0$ .

151 **Remark 3.2.** We will describe a best-of-both-worlds type algorithm, that attains optimal guarantees  
 152 both under stochastic and adversarial inputs, without knowledge of the specific setting in which the  
 153 algorithm operates. It should be noted that  $\rho_{\text{Adv}}$  and  $\rho_{\text{Stoc}}$  are *not* known by the algorithm. While  
 154 the algorithm could potentially efficiently estimate  $\rho_{\text{Stoc}}$  in stochastic settings, as shown in [19],  
 155 acquiring knowledge of  $\rho_{\text{Adv}}$  in the adversarial setting would necessitate information about future  
 156 inputs. This requirement is generally unfeasible for most instances of interest.

## 157 4 On Best-Of-Both-Worlds Guarantees

158 We employ the expression *best-of-both-worlds* as defined in [14] for the case of online allocation  
 159 problems with resource-consumption constraints. In this context, we expect different types of  
 160 guarantees depending on the input model being considered.

161 When inputs are stochastic, a best-of-both-worlds algorithm should guarantee that, given failure  
 162 probability  $\delta > 0$ , with probability at least  $1 - \delta$

$$\max(\text{Opt}_{\text{Stoc}} - \text{Rew}(T), V(T)) = \tilde{O}(\sqrt{T}).$$

163 The dependency on  $T$  is optimal since, in the worst case, it is optimal even without constraints [7].

164 In adversarial settings, a best-of-both-worlds algorithm should guarantee that, with probability at  
 165 least  $1 - \delta$ ,

$$\max(\text{Opt}_{\text{Adv}} - \alpha \text{Rew}(T), V(T)) = \tilde{O}(\sqrt{T}),$$

166 where  $\alpha > 1$  is the *competitive ratio*. In the BwK scenario with only resource-consumption constraints,  
 167 the optimal competitive ratio attainable is  $\alpha = 1/\rho_{\text{Adv}}$ . In that setting,  $\rho_{\text{Adv}}$  denotes the per-iteration  
 168 budget, which we can assume is equal for each resource without loss of generality. In our set-up,  
 169 considering arbitrary and potentially negative constraints, we will present an algorithm for which the  
 170 above holds for  $\alpha := 1 + 1/\rho_{\text{Adv}}$ . The following result shows that this competitive ratio is optimal.  
 171 In particular, we show that it is not possible to obtain cumulative constraint violations of order  $o(T)$   
 172 and competitive ratio strictly less than  $1 + 1/\rho_{\text{Adv}}$  (omitted proofs can be found in the Appendix).

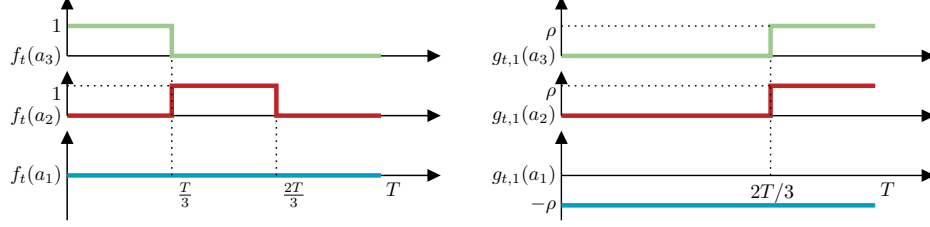


Figure 1: Reward and costs of each arm of the instance employed in Example 5.2.

173 **Theorem 4.1.** [Lower bound adversarial setting] Consider the family of all adversarial instances  
 174 with  $\mathcal{X} = \{a_1, a_2\}$ , each characterized by a parameter  $\rho_{\text{Adv}}$  and optimal reward  $\text{Opt}_{\text{Adv}}$ . Then,  
 175 no algorithm can achieve, on all instances, sublinear cumulative violations  $\mathbb{E}[V(T)] = o(T)$  and  
 176  $\text{Opt}_{\text{Adv}} \mathbb{E}[\text{Rew}] > 1 + 1/\rho_{\text{Adv}}$ .

## 177 5 Lagrangian Framework

178 Given the reward function  $f : \mathcal{X} \rightarrow [0, 1]$  and  
 179 the costs functions  $g : \mathcal{X} \rightarrow [-1, 1]^m$  we define  
 180 the Lagrangian  $\mathcal{L}_{f,g} : \mathcal{X} \times \mathbb{R}_+^m \rightarrow \mathbb{R}$  as:

$$\mathcal{L}_{f,g}(x, \lambda) := f(x) - \langle \lambda, g(x) \rangle.$$

181 We will consider a modular primal-dual ap-  
 182 proach that employs a *primal* algorithm  $\text{Alg}_P$ ,  
 183 producing primal decisions  $x_t$ , and a *dual* algo-  
 184 rithm  $\text{Alg}_D$  that produces dual decisions  $\lambda_t$  for  
 185 all  $t$ . We assume that  $\text{Alg}_P$  and  $\text{Alg}_D$  produce  
 186 their decisions in order to maximize their utili-  
 187 ties  $u_t^P$  and  $u_t^D$ , respectively. We define  $u_t^P : x \mapsto$   
 188  $\mathcal{L}_{f_t, g_t}(x, \lambda_t)$  and  $u_t^D : \lambda \mapsto -\mathcal{L}_{f_t, g_t}(x_t, \lambda)$ .  
 189 The regret of the primal algorithm  $\text{Alg}_P$  on any  
 190 subset  $I \subseteq \llbracket T \rrbracket$  is defined as:

$$R_I^P(\mathcal{X}) := \sup_{x \in \mathcal{X}} \sum_{t \in I} [u_t^P(x) - u_t^P(x_t)].$$

191 The regret of the dual algorithm  $\text{Alg}_D$  is defined similarly for any bounded subset  $\mathcal{D} \subseteq \mathbb{R}_+$ :

$$R_I^D(\mathcal{D}) := \sup_{\lambda \in \mathcal{D}} \sum_{t \in I} [u_t^D(\lambda) - u_t^D(\lambda_t)].$$

192 For ease of notation we write  $R_T^P(\mathcal{X})$  and  $R_T^D(\mathcal{D})$  when  $I = \llbracket T \rrbracket$ , instead of  $R_{\llbracket T \rrbracket}^P(\mathcal{X})$  and  $R_{\llbracket T \rrbracket}^D(\mathcal{D})$ .

193 The interaction of  $\text{Alg}_P$  and  $\text{Alg}_D$  with the environment is reported in Algorithm 1. Note that the  
 194 feedback of  $\text{Alg}_P$  is forced to be bandit by the fact that we do not have counterfactual information of  
 195  $f_t$  and  $g_t$ , however  $\text{Alg}_D$  receives full feedback by design.

196 **Remark 5.1** (The Challenges of the Adversarial Setting). In the stochastic setting, it is not required  
 197 adaptive regret minimization, see *e.g.*, [40], as it is possible to analyze directly the expected zero-sum  
 198 game between  $\text{Alg}_P$  and  $\text{Alg}_D$ . However, in the adversarial setting, the algorithms  $\text{Alg}_P$  and  $\text{Alg}_D$   
 199 face a different zero-sum game at each time  $t$ . Indeed, since  $f_t$  and  $g_t$  are adversarial, the zero-sum  
 200 game with payoffs  $\mathcal{L}_{f_t, g_t}(\cdot, \cdot)$  is only seen at time  $t$ . This is in contrast to what happens in the  
 201 stochastic setting in which the zero-sum game  $\mathcal{L}_{\bar{f}, \bar{g}}(\cdot, \cdot)$  at each time  $t$  is the same for all time  $t$ .

### 202 5.1 No-Regret is Not Enough!

203 Typically, Lagrangian frameworks for constrained bandit problems are solved by instantiating  $\text{Alg}_P$   
 204 and  $\text{Alg}_D$  with two regret minimizers, which are algorithms guaranteeing  $R_T^P(\mathcal{X}), R_T^D(\mathcal{D}) = o(T)$ ,  
 205 respectively [30, 17]. The dual regret minimizer is usually instantiated with  $\mathcal{D} := [0, M]^m$ , for some  
 206 constant  $M > 0$ . Ensuring that  $\mathcal{D}$  is bounded is crucial to control the magnitude of primal utilities

---

#### Algorithm 1 Primal-Dual Algorithm

---

- 1: **Input:**  $\text{Alg}_P$  and  $\text{Alg}_D$ .
  - 2: **for**  $t = 1, 2, \dots, T$  **do**
  - 3:   **Primal decision:**  $x_t \leftarrow \text{Alg}_P$
  - 4:   **Dual decision:**  $\lambda_t \leftarrow \text{Alg}_D$
  - 5:   **Observe:**  $f_t(x_t)$  and  $g_t(x_t)$
  - 6:   **Primal update:** feed  $u_t^P(x_t)$  to  $\text{Alg}_P$ ,  
       where
  - 7:      $u_t^P(x_t) \leftarrow f_t(x_t) - \langle \lambda_t, g_t(x_t) \rangle$
  - 8:   **Dual update:**  
       Feed  $u_t^D : \lambda \mapsto -f_t(x_t) + \langle \lambda, c_t(x_t) \rangle$  to  
        $\text{Alg}_D$
  - 9: **end for**
-

207  $u_t^p(\cdot)$ , whose scale influences the magnitude of the primal regret. In the following example, we show  
 208 that we cannot rely solely on arguments based on the *black-box* no-regret property of  $\text{Alg}_P$  and  $\text{Alg}_D$   
 209 and hence we need stronger guarantees than simple no-regret.

210 **Example 5.2.** We have one constraint, i.e.,  $m = 1$  and the set  $\mathcal{X} = \{a_1, a_2, a_3\}$  is a discrete set of 3  
 211 actions. The rewards of  $a_1$  is always 0, i.e.,  $f_t(a_1) = 0$  for all  $t \in \llbracket T \rrbracket$ , while its cost is always  $-\rho$ ,  
 212 i.e.,  $g_{t,1}(a_1) = -\rho$  for all  $t \in t$ . The rewards for  $a_2$  and  $a_3$  are defined as follows: for  $t \in \llbracket T/3 \rrbracket$   
 213 we have  $f_t(a_2) = 0$  while  $f_t(a_3) = 1$ . On the other hand, for  $t \in \llbracket T/3, 2T/3 \rrbracket$  we have  $f_t(a_2) = 1$   
 214 while  $f_t(a_3) = 0$ . Finally  $f_t(a_2) = f_t(a_3) = 0$  for all  $t \in \llbracket 2T/3, T \rrbracket$ . The costs for  $a_2$  and  $a_3$  are  
 215 defined as follows: for  $t \in \llbracket 2T/3 \rrbracket$  we have  $g_{t,1}(a_2) = g_{t,1}(a_3) = 0$ , while  $g_{t,1}(a_2) = g_{t,1}(a_3) = 1$   
 216 for all  $t \in \llbracket 2T/3, T \rrbracket$ . The instance is depicted in Figure 1.

217 **Proposition 5.3.** Consider the instance of Example 5.2. Even if  $\text{Alg}_P$  and  $\text{Alg}_D$  suffer regret less  
 218 than or equal then zero, the primal-dual framework fails to achieve sublinear constraint violations.

219 Intuitively, the reason for which a standard primal-dual framework fails in Example 5.2 is that the  
 220 primal regret minimizer can accumulate enough negative regret in the first two phases to “absorb”  
 221 large regret suffered in the third phase. This “laziness” of  $\text{Alg}_P$  allows it to play actions in the  
 222 last phase for which it incurs linear violations of the constraint. For more details see the proof of  
 223 Proposition 5.3 in Appendix A. One could solve the problem employing the *recovery technique*  
 224 proposed in [19], which prescribes to minimize the violations at a prescribed time. However, selecting  
 225 the right time to start the recovery phase crucially requires knowledge of the Slater’s parameter,  
 226 which is not available in our setting. The only approach which does not require knowledge of Slater’s  
 227 parameter is the one proposed in [18] for the case of *return-on-investment* constraints, whose core  
 228 idea we describe in the next section.

229 **Remark 5.4.** We remark that it is not possible to prove that any choice of  $\text{Alg}_P$  and  $\text{Alg}_D$  satisfying  
 230 the no-regret property fails in our setting. Indeed, we will end up choosing  $\text{Alg}_P$  and  $\text{Alg}_D$  algorithms  
 231 that have a *stronger* no-regret property (and hence are also no-regret). Proposition 5.3 shows that  
 232 our arguments and algorithms must necessarily rely on a stronger version of regret, specifically  
 233 *no-adaptive regret*.

## 234 5.2 No-Adaptive Regret

235 The reason why generic regret minimizers fail to give satisfactory result on the instance described in  
 236 Example 5.2 is that they fail to adapt to the changing environment, even if the regret of the primal  
 237 is zero on the entire horizon  $\llbracket T \rrbracket$ , it fails to “adapt” in the final rounds  $\llbracket 2T/3, T \rrbracket$ . Indeed, in these  
 238 last rounds, if the primal algorithm’s objective is guaranteeing sublinear regret over  $\llbracket T \rrbracket$ , it is not  
 239 required to update its decision, since it accumulated large negative regret of  $-2T/3$  regret in the  
 240 initial rounds  $\llbracket 2T/3 \rrbracket$ . Therefore, standard no-regret guarantees are not enough.

241 A stronger requirement for the primal and dual algorithm is being *weakly adaptive* [29], that is,  
 242 guaranteeing that in high probability  $\sup_{I=\llbracket t_1, t_2 \rrbracket} R_I^{P,D} = o(T)$ . Intuitively, this requirement would  
 243 force  $\text{Alg}_P$  to change its action during the last phase of Example 5.2. This idea was first proposed in  
 244 [18] for the specific case of a learner with one budget and one return-on-investments constraints. In  
 245 the following section, we show how such approach can be extended to the case of general constraints.

## 246 6 Self-Bounding Lemma

247 One crucial difference with the previous literature is that the feasibility parameter is not known a  
 248 priori, and thus we cannot directly bound the range of the Lagrange multipliers as in BwK. At a high  
 249 level we want that, regardless of the choices of  $f_t$  and  $g_t$ , the  $\ell_1$  norm of the Lagrange multipliers  
 250 is bounded by a quantity that depends on the (unknown) parameters of the instance. However, for  
 251 this to hold we need that the primal algorithm  $\text{Alg}_P$  is (almost) scale free, i.e., that its regret scale  
 252 quadratically in the unknown range of its reward function.<sup>2</sup> Formally:

253 **Definition 6.1.** For any  $c \geq 1$ , we say that  $\text{Alg}_P$  is a  $c$ -scale-free and weakly-adaptive regret  
 254 minimizer if, for any subset of rounds  $I = \llbracket t_1, t_2 \rrbracket \subseteq \llbracket T \rrbracket$ , with probability at least  $1 - \delta$  it holds that

$$R_I^P(\mathcal{X}) \leq L^c \cdot \overline{R}_{T,\delta}^P(\mathcal{X}),$$

<sup>2</sup>Usually we say that an algorithm is scale-free [37] if its regret scales linearly in the (unknown) range of its rewards, i.e., 1-scale-free with our definition.

255 where the maximum module of the primal utilities is  $\sup_{t \in \llbracket T \rrbracket, x \in \mathcal{X}} |u_t^p(x)| =: L$ , and  $\overline{R}_{T,\delta}^p(\mathcal{X})$   
 256 depends only on  $T$ ,  $\delta$  and  $\mathcal{X}$ , and is non-decreasing in the length of the time horizon  $T$ .

257 Now, we show that *online gradient descent* (OGD) [46] with a carefully defined learning rate yields  
 258 the required self-bounding property both in the stochastic and adversarial setting.

259 **Lemma 6.2** (Self-bounding lemma). *Let  $\eta_{\text{OGD}} := (800 \cdot m \cdot \max\{\overline{R}_{T,\delta}^p(\mathcal{X}), E_{T,\delta}\})^{-1}$ , then if  
 260  $\text{Alg}_{\mathcal{D}}$  is OGD on the set  $\mathcal{D} = \mathbb{R}_{\geq 0}^m$ , and the primal algorithm  $\text{Alg}_{\mathcal{P}}$  is 2-scale-free and has a  
 261 high-probability weakly adaptive regret bound  $\overline{R}_{T,\delta}^p(\mathcal{X})$ , then with probability at least  $1 - \delta$ :*

$$\max_{t \in \llbracket T \rrbracket} \|\lambda_t\|_1 \leq \frac{13m}{\rho},$$

262 where  $\rho = \rho_{\text{Adv}}$  or  $\rho = \rho_{\text{Stoc}}$  depending on the setting and  $E_{T,\delta} := \sqrt{16T \log(2T/\delta)}$ .

263 We remark that the self-bounding lemma shows that, if we take OGD with a carefully defined learning  
 264 rate  $\eta_{\text{OGD}} = \tilde{O}((m \max\{\overline{R}_{T,\delta}^p(\mathcal{X}), \sqrt{T}\})^{-1})$  as  $\text{Alg}_{\mathcal{P}}$ , then the  $\ell_1$ -norm of the variables  $\lambda_t$  is  
 265 automatically bounded by the reciprocal of the feasibility parameter, even if the feasibility parameter  
 266 is unknown to the learner. This is the central result that allows us to build algorithms that work  
 267 without knowing Slater's parameter. We observe that:

268 **Remark 6.3.** Even in the simplest instances of bandit problems one has  $\overline{R}_{T,\delta}^p(\mathcal{X}) = \tilde{\Omega}(\sqrt{T})$  and,  
 269 therefore, we can assume that  $\eta_{\text{OGD}} = \tilde{O}((m \overline{R}_{T,\delta}^p(\mathcal{X}))^{-1})$ .

270 **Remark 6.4.** We will work with 2-scale-free algorithms, which suffice to obtain the desired guaran-  
 271 tees for our framework. We observe that scale-free algorithms would yield a tighter bound of  $1/\rho$  in  
 272 the Theorems 7.2 and 7.3 and a simpler analysis of Lemma 6.2. However, scale-free algorithm are  
 273 much more difficult to find and this would limit the extent to which our framework can be applied.  
 274 On the other hand, 2-scale-free algorithm seems to be more abundant (see, *e.g.*, Section 8). Indeed,  
 275 as we show in Section 8, it is usually the case that setting the learning rate independent on the scale  
 276 of the rewards provides 2-scale-freeness. We leave such characterization to future research.

## 277 7 General Guarantees

278 First, we exploit Lemma 6.2 to bound the total violations of the framework.

279 **Theorem 7.1.** *Let  $\text{Alg}_{\mathcal{D}}$  be OGD with learning rate  $\eta$  as in Lemma 6.2, and let  $\text{Alg}_{\mathcal{P}}$  any 2-  
 280 scale-free algorithm with no-adaptive regret. Then, with probability at least  $1 - \delta$ , it holds that  
 281  $V_T = \tilde{O}\left(\frac{m^2}{\rho} \overline{R}_{T,\delta}^p(\mathcal{X})\right)$ , where  $\rho = \rho_{\text{Adv}}$  in the adversarial setting and  $\rho = \rho_{\text{Stoc}}$  in the stochastic.*

282 Moreover, the proof of Theorem 7.1 can be easily adapted to show that the violations of any constraint  
 283  $i \in \llbracket m \rrbracket$  is bounded on any interval  $\llbracket t \rrbracket$  with  $t \in \llbracket T \rrbracket$ .

284 Now, we prove that the framework, with high probability, yields optimal guarantees in both stochastic  
 285 and adversarial settings. We start with the adversarial setting, for which the following result holds.

286 **Theorem 7.2.** *If  $\text{Alg}_{\mathcal{D}}$  is OGD with learning rate  $\eta_{\text{OGD}}$  and domain  $\mathcal{D} := \mathbb{R}_{\geq 0}^m$ , and  $\text{Alg}_{\mathcal{P}}$  is 2-scale-  
 287 free, then, in the adversarial setting, with high probability:*

$$\text{Rew} \geq \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} - \tilde{O}\left(\left(\frac{m}{\rho_{\text{Adv}}}\right)^2 \overline{R}_{T,\delta}^p(\mathcal{X})\right).$$

288 On the other hand, for the stochastic setting we can prove the following result:

289 **Theorem 7.3.** *If  $\text{Alg}_{\mathcal{D}}$  is OGD with learning rate  $\eta_{\text{OGD}}$  and domain  $\mathcal{D} := \mathbb{R}_{\geq 0}^m$ , and  $\text{Alg}_{\mathcal{P}}$  is 2-scale-  
 290 free, then in the stochastic setting, in high probability:*

$$\text{Rew} \geq \text{Opt}_{\text{Stoc}} - \tilde{O}\left(\left(\frac{m}{\rho_{\text{Stoc}}}\right)^2 \overline{R}_{T,\delta}^p(\mathcal{X})\right).$$

291 **Remark 7.4.** Any algorithm with vanishing constraints violations can be employed to handle also  
 292 BwK constraints. In such setting, the learner has resource-consumption constraints with *hard stopping*  
 293 (*i.e.*, once the budget for a resource is fully depleted the learner must play the void action until the  
 294 end of time horizon). This does not yield any fundamental complication for our framework. Indeed,  
 295 we could introduce an initial phase of  $o(T)$  rounds in which the algorithm collects the extra budget  
 296 needed to cover potential violations, before starting the primal-dual procedure.

## 297 8 Applications

298 In this section, we show how our framework can be instantiated to handle scenarios such as bandits  
 299 with general constraints, as well as contextual bandits with constraints (*i.e.*,  $\text{CBwLC}$ ). Thanks to the  
 300 modularity of the results derived in the previous sections, we only need to provide an algorithm  $\text{Alg}_P$   
 301 which is 2-scale-free and weakly adaptive for a desired action space  $\mathcal{X}$  and rewards  $u_t^P$ .

### 302 8.1 Bandits with General Constraints

303 In this setting, the action space is  $\mathcal{X} = \llbracket K \rrbracket$ . [18] showed that the  $\text{EXP3-SIX}$  algorithm introduced  
 304 by [36] can be used as  $\text{Alg}_P$ , since it guarantees sublinear weakly adaptive regret in high probability,  
 305 and it is 2-scale-free.

306 **Theorem 8.1** (Theorem 8.1 of [18]). *EXP3-SIX instantiated with suitable parameters guarantees*  
 307 *that, with probability at least  $1 - \delta$  that  $\sup_{I=\llbracket t_1, t_2 \rrbracket} R_I^P(\mathcal{X}) = O\left(\sqrt{KT} \log(KT\delta^{-1})\right)$ .*

308 Thus, by applying Theorem 7.1 on the violations, and Theorem 7.2 and Theorem 7.3 on the adversarial  
 309 and stochastic reward guarantees respectively, we get the following result:

310 **Corollary 8.2.** *Consider a multi armed bandit problem with constraints. There exists an algorithm*  
 311 *that w.h.p. guarantees, in the adversarial setting, violations at most  $\tilde{O}\left(\frac{m^2}{\rho_{\text{Adv}}}\sqrt{KT}\right)$  and  $R_{\text{ew}} \geq$*   
 312  *$\frac{\rho_{\text{Adv}}}{1+\rho_{\text{Adv}}} \text{opt}_{\text{Adv}} - \tilde{O}\left(\frac{m^2}{\rho_{\text{Adv}}^2}\sqrt{KT}\right)$ , while, in the stochastic setting, it guarantees violations at most*  
 313  *$\tilde{O}\left(\frac{m^2}{\rho_{\text{Stoc}}}\sqrt{KT}\right)$  and reward at least  $R_{\text{ew}} \geq \text{opt}_{\text{Stoc}} - \tilde{O}\left(\frac{m^2}{\rho_{\text{Stoc}}^2}\sqrt{KT}\right)$ .*

### 314 8.2 Contextual Bandits with Constraints

315 Following [41], we apply our general framework to contextual bandits with regression oracles. In  
 316 this setting, the decision maker observes a context  $z_t \in \mathcal{Z}$  from some context set  $\mathcal{Z}$ , where  $z_t$  is  
 317 possibly chosen by an adversary. Then, the decision maker picks its decision  $a_t$  from an action set  $\mathcal{A}$ .  
 318 Then, the reward is computed as a function of the context and the action, *i.e.*,  $f_t : \mathcal{Z} \times \mathcal{A} \rightarrow [0, 1]$ ,  
 319 and similarly for the constraints  $\mathbf{g}_t : \mathcal{Z} \times \mathcal{A} \rightarrow [-1, 1]^m$ . At each  $t$ ,  $f_t$  and  $\mathbf{g}_t$  are drawn from  
 320 some distribution. More precisely, there exist a class  $\mathcal{F}$  of functions and  $\bar{f}, \bar{g}_i \in \mathcal{F}$  such that for all  
 321  $(z, a) \in \mathcal{Z} \times \mathcal{A}$  it holds that  $\mathbb{E}[f_t(z, a)|z, a] = \bar{f}(z, a)$  and  $\mathbb{E}[g_{t,i}(z, a)|z, a] = \bar{g}_i(z, a)$  for  $i \in \llbracket m \rrbracket$ .

322 We slightly modify the primal-dual algorithm to handle contexts. In particular,  $\text{Alg}_P$  gets to observe  
 323 a context  $z_t$  before deciding their action. Formally, we can use the machinery introduced in Section 3  
 324 by taking  $\mathcal{X}$  as the set of deterministic policies  $\Pi := \{\pi : \mathcal{Z} \rightarrow \mathcal{A}\}$ . Then,  $u_t^P(\pi) = f_t(z_t, \pi(z_t)) -$   
 325  $\langle \boldsymbol{\lambda}_t, \mathbf{g}_t(z_t, \pi(z_t)) \rangle$ , and the action  $a_t$  is computed through  $\pi_t$  returned by the primal algorithm.  
 326 Although this choice transforms the contextual framework into an application of the framework  
 327 introduced in Section 3, in practical terms, it is simpler to think of  $a_t$  as the direct output of  $\text{Alg}_P$   
 328 upon observing the context  $z_t$ . The extended primal-dual framework is sketched in Algorithm 2.

329 We assume to have  $m + 1$  online regression oracles  $(\mathcal{O}_f, \mathcal{O}_1, \dots, \mathcal{O}_m)$  for the functions  $\bar{f}$  and  
 330  $\bar{g}_1, \dots, \bar{g}_m$ , respectively. The regression oracle  $\mathcal{O}_f$  produces, at each  $t$ , a regressor  $\hat{f}_t \in \mathcal{F}$  that tries  
 331 to approximate the *true* regressor  $\bar{f}$ . Then, the oracle is feed with a new data point, comprised of a  
 332 context  $z_t \in \mathcal{Z}$  and an action  $a_t \in \mathcal{A}$ , and the performance of the regressor is evaluated on the basis  
 333 of its prediction for the tuple  $(z_t, a_t)$ . The online regression oracle  $\mathcal{O}_f$  is updated with the labeled  
 334 data point  $(z_t, a_t, f_t(z_t, a_t))$ . Overall, its performance is measured by its cumulative  $\ell_2$ -error:

$$\text{Err}(\mathcal{O}_f) := \sum_{t \in \llbracket T \rrbracket} \left( \hat{f}_t(z_t, a_t) - \bar{f}(z_t, a_t) \right)^2.$$

335 Each online regression oracle  $(\mathcal{O}_i)_{i \in \llbracket m \rrbracket}$  works analogously, and its performance is measured by  
 336  $\text{Err}(\mathcal{O}_i) := \sum_{t \in \llbracket T \rrbracket} (\hat{g}_t(z_t, a_t) - \bar{g}_i(z_t, a_t))^2$ .

337 By combining the online regression oracles  $\mathcal{O}_f$  and  $\{\mathcal{O}_i\}_{i \in \llbracket m \rrbracket}$  we can build an online regression  
 338 oracle  $\mathcal{O}_{\mathcal{L}}$  for the Lagrangian which outputs regressors  $\hat{\mathcal{L}}_t : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}$  defined as:

$$\hat{\mathcal{L}}_t(z, a) = \mathcal{L}_{\hat{f}_t, \hat{\mathbf{g}}_t}((z, a), \boldsymbol{\lambda}_t) = \hat{f}_t((z, a)) - \langle \boldsymbol{\lambda}_t, \hat{\mathbf{g}}_t(z, a) \rangle,$$



---

**Algorithm 2** Primal-Dual Algorithm for Contextual Bandits
 

---

1: **Input:**  $\text{Alg}_P$  and  $\text{Alg}_D$ .  
 2: **for**  $t = 1, 2, \dots, T$  **do**  
 3:   Observe context  $z_t$   
 4:   **Dual decision:**  $\lambda_t \leftarrow \text{Alg}_D$   
 5:   **Primal decision:**  
 6:      $a_t \leftarrow \text{Alg}_P(z_t, \lambda_t)$   
 7:   **Observe:**  $f_t(z_t, a_t)$  and  $g_t(z_t, a_t)$   
 8:   **Primal update:** feed  $u_t^P(a_t)$  to  $\text{Alg}_P$ ,  
   where  
 9:      $u_t^P(a_t) = f_t(z_t, a_t) - \langle \lambda_t, g_t(z_t, a_t) \rangle$   
 10:   **Dual update:** feed  $u_t^D$  to  $\text{Alg}_D$ ,  
   where  
 10:      $u_t^D(\lambda) = f_t(z_t, a_t) + \langle \lambda, c_t(z_t, a_t) \rangle$   
 11: **end for**

---



---

**Algorithm 3** Primal Algorithm for Contextual Bandits
 

---

1: **Input:** Learning rate  $\eta_P$   
 2: **Get regressors from online regression oracles:**  
 3:    $\hat{f}_t \leftarrow \mathcal{O}_f$ , and  $\hat{g}_{t,i} \leftarrow \mathcal{O}_i$  for all  $i \in \llbracket m \rrbracket$   
 4: Observe context  $z_t$  and dual variable  $\lambda_t$   
 5: For all  $a \in \mathcal{A}$  compute  $\hat{\mathcal{L}}_t(a) := \mathcal{L}_{\hat{f}_t, \hat{g}_t}((z_t, a), \lambda_t)$   
 6: Compute  $\xi_t \in \Delta(\mathcal{A})$  as:  
   
$$\xi_t(a) = \left( \mu_t + \eta_P \left( \max_{a'} \hat{\mathcal{L}}_t(a') - \hat{\mathcal{L}}_t(a) \right) \right)^{-1}$$
  
    $\triangleright \mu_t$  is such that  $\xi_t \in \Delta(\mathcal{A})$   
 7: Sample  $a_t \sim \xi_t$  and return it.  
 8: **Update online regression oracles:**  
 9:   Feed  $(z_t, a_t, f_t(z_t, a_t))$  to  $\mathcal{O}_f$   
 10:   Feed  $(z_t, a_t, g_{t,i}(z_t, a_t))$  to  $\mathcal{O}_i \forall i \in \llbracket m \rrbracket$

---

339 while we define  $\tilde{\mathcal{L}}(z, a) := \mathcal{L}_{\tilde{f}, \tilde{g}}((z, a), \lambda_t)$ . The  $\ell_2$ -error of  $\mathcal{O}_{\mathcal{L}}$  can be bounded via the following  
 340 extension of [40, Theorem 16].

341 **Lemma 8.3.** *The error of  $\mathcal{O}_{\mathcal{L}}$  can be bounded as*

$$\text{Err}(\mathcal{O}_{\mathcal{L}}) \leq 2\text{Err}(\mathcal{O}_f) + 2 \left( \sup_{t \in \llbracket T \rrbracket} \|\lambda_t\|_1 \right)^2 \sum_{i \in \llbracket m \rrbracket} \text{Err}(\mathcal{O}_i).$$

342 The fundamental idea of [25] is to reduce (unconstrained) contextual bandit problems to online linear  
 343 regression. Recently, this ideas was extended in [41, 27] in order to design a primal algorithm  $\text{Alg}_P$   
 344 capable of handling stochastic contextual bandits with constraints (see Algorithm 3).

345 To apply Algorithm 3 to our framework we need to find an algorithm  $\text{Alg}_P$  which is 2-scale-free and  
 346 weakly adaptive with high probability. We extend the result [25] to prove that their reduction actually  
 347 satisfies the required guarantees.

348 **Lemma 8.4.** *Assume that  $\max\{\text{Err}(\mathcal{O}_f), \text{Err}(\mathcal{O}_i)\} \leq \overline{\text{Err}}$ . Then, we have that Algorithm 3 with  
 349  $\eta_P := \sqrt{KT}$  guarantees that  $\sup_{I=\llbracket t_1, t_2 \rrbracket} R_I^P(\Pi) = \tilde{O}\left(m \cdot \overline{\text{Err}} \cdot L^2 \cdot \sqrt{KT}\right)$  with high probability,  
 350 where  $L := \sup_{t \in \llbracket T \rrbracket, \pi \in \Pi} |u_t^P(\pi)|$ .*

351 Equipped with a 2-scale free algorithm that suffers no adaptive regret with high probability, we  
 352 can combine  $\text{Alg}_P$  with the results of Theorems 7.1 to 7.3 to prove the first optimal guarantees for  
 353 CBwLC with adversarial contexts.

354 **Corollary 8.5.** *Consider a functional class  $\mathcal{F}$  and an online regression oracle that guarantees  $\ell_2$ -  
 355 error  $\overline{\text{Err}}$ . There exists an algorithm that w.h.p. guarantees violations at most  $\tilde{O}\left(\frac{m^3}{\rho_{\text{Adv}}} \overline{\text{Err}} \sqrt{KT}\right)$   
 356 and reward at least  $\text{Rew} \geq \frac{\rho_{\text{Adv}}}{1+\rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} - \tilde{O}\left(\overline{\text{Err}} \frac{m^3}{\rho_{\text{Adv}}^2} \sqrt{KT}\right)$  in the adversarial setting, while  
 357 it guarantees violations at most  $\tilde{O}\left(\frac{m^3}{\rho_{\text{Stoc}}} \overline{\text{Err}} \sqrt{KT}\right)$  and reward at least  $\text{Rew} \geq \text{Opt}_{\text{Stoc}} -$   
 358  $\tilde{O}\left(\overline{\text{Err}} \frac{m^3}{\rho_{\text{Stoc}}^2} \sqrt{KT}\right)$  in the stochastic setting.*

359 [25] includes many examples of functional classes  $\mathcal{F}$  that have good online regression oracles,  
 360 meaning that their error is subpolynomial in the time horizon  $T$ . We report here some notable  
 361 mentions for completeness.

362 If  $\mathcal{F}$  is a finite set of functions we have that  $\overline{\text{Err}} = O(\log |\mathcal{F}|)$ , which comes from using as regression  
 363 oracles the Vovk forecaster [42]. Another important examples is the case in which  $\mathcal{F}$  is the class of  
 364 linear functions, i.e.,  $\mathcal{F} = \{h(z, a) = \langle z_a, \theta \rangle : \theta \in \mathbb{R}^d, \|\theta\|_2 \leq 1\}$ , i.e., each actions  $a$  is associated  
 365 with a known feature vector  $z_a \in \mathbb{R}^d$  which generates the reward/costs trough a unknown parameter  
 366  $\theta$  that characterize the linear function. Here, there exists a online regression oracle which provides  
 367  $\ell_2$ -error  $\overline{\text{Err}} = O(d \log(T/d))$  [8].

368 **References**

- 369 [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear  
370 stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- 371 [2] Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire.  
372 Taming the monster: A fast and simple algorithm for contextual bandits. In *International  
373 Conference on Machine Learning*, pages 1638–1646. PMLR, 2014.
- 374 [3] Shipra Agrawal and Nikhil Devanur. Linear contextual bandits with knapsacks. *Advances in  
375 Neural Information Processing Systems*, 29, 2016.
- 376 [4] Shipra Agrawal and Nikhil R Devanur. Bandits with concave rewards and convex knapsacks. In  
377 *Proceedings of the fifteenth ACM conference on Economics and computation*, pages 989–1006.  
378 ACM, 2014.
- 379 [5] Shipra Agrawal and Nikhil R Devanur. Bandits with global convex constraints and objective.  
380 *Operations Research*, 67(5):1486–1502, 2019.
- 381 [6] Shipra Agrawal, Nikhil R Devanur, and Lihong Li. An efficient algorithm for contextual bandits  
382 with knapsacks, and an extension to concave objectives. In *29th Annual Conference on Learning  
383 Theory (COLT)*, 2016.
- 384 [7] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic  
385 multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- 386 [8] Katy S Azoury and Manfred K Warmuth. Relative loss bounds for on-line density estimation  
387 with the exponential family of distributions. *Machine learning*, 43:211–246, 2001.
- 388 [9] Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. Bandits with knap-  
389 sacks. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science, FOCS  
390 2013*, pages 207–216. IEEE, 2013.
- 391 [10] Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. Bandits with knap-  
392 sacks. *J. ACM*, 65(3), 2018.
- 393 [11] Ashwinkumar Badanidiyuru, John Langford, and Aleksandrs Slivkins. Resourceful contextual  
394 bandits. In *Conference on Learning Theory*, pages 1109–1134. PMLR, 2014.
- 395 [12] Santiago Balseiro, Christian Kroer, and Rachitesh Kumar. Online resource allocation under  
396 horizon uncertainty. *SIGMETRICS Perform. Eval. Rev.*, 51(1):63–64, 2023.
- 397 [13] Santiago R Balseiro and Yonatan Gur. Learning in repeated auctions with budgets: Regret  
398 minimization and equilibrium. *Management Science*, 65(9):3952–3968, 2019.
- 399 [14] Santiago R Balseiro, Haihao Lu, and Vahab Mirrokni. The best of many worlds: Dual mirror  
400 descent for online allocation problems. *Operations Research*, 2022.
- 401 [15] Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. Bandits with  
402 replenishable knapsacks: the best of both worlds. *arXiv preprint arXiv:2306.08470*, 2023.
- 403 [16] Alberto Bietti, Alekh Agarwal, and John Langford. A contextual bandit bake-off. *The Journal  
404 of Machine Learning Research*, 22(1):5928–5976, 2021.
- 405 [17] Matteo Castiglioni, Andrea Celli, and Christian Kroer. Online learning with knapsacks: the best  
406 of both worlds. In *International Conference on Machine Learning*, pages 2767–2783. PMLR,  
407 2022.
- 408 [18] Matteo Castiglioni, Andrea Celli, and Christian Kroer. Online bidding in repeated non-truthful  
409 auctions under budget and ROI constraints. *arXiv preprint arXiv:2302.01203*, 2023.
- 410 [19] Matteo Castiglioni, Andrea Celli, Alberto Marchesi, Giulia Romano, and Nicola Gatti. A  
411 unifying framework for online optimization with long-term constraints. In *Advances in Neural  
412 Information Processing Systems*, volume 35, pages 33589–33602, 2022.

- 413 [20] Andrea Celli, Matteo Castiglioni, and Christian Kroer. Best of many worlds guarantees for  
414 online learning with knapsacks. *arXiv preprint arXiv:2202.13710*, 2023.
- 415 [21] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff  
416 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence  
417 and Statistics*, pages 208–214, 2011.
- 418 [22] Zhe Feng, Swati Padmanabhan, and Di Wang. Online bidding algorithms for return-on-spend  
419 constrained advertisers. In *Proceedings of the ACM Web Conference 2023*, page 3550–3560,  
420 2023.
- 421 [23] Giannis Fikioris and Éva Tardos. Approximately stationary bandits with knapsacks. In *Proceed-  
422 ings of Thirty Sixth Conference on Learning Theory*, volume 195, pages 3758–3782, 12–15 Jul  
423 2023.
- 424 [24] Dylan Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert Schapire. Practical  
425 contextual bandits with regression oracles. In *International Conference on Machine Learning*,  
426 pages 1539–1548. PMLR, 2018.
- 427 [25] Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits  
428 with regression oracles. In *International Conference on Machine Learning*, pages 3199–3210.  
429 PMLR, 2020.
- 430 [26] Jason Gaitonde, Yingkai Li, Bar Light, Brendan Lucier, and Aleksandrs Slivkins. Budget  
431 pacing in repeated auctions: Regret and efficiency without convergence. In *14th Innovations in  
432 Theoretical Computer Science Conference (ITCS 2023)*, volume 251, page 52, 2023.
- 433 [27] Yuxuan Han, Jialin Zeng, Yang Wang, Yang Xiang, and Jiheng Zhang. Optimal contextual  
434 bandits with knapsacks under realizability via regression oracles. In *International Conference  
435 on Artificial Intelligence and Statistics*, pages 5011–5035. PMLR, 2023.
- 436 [28] Elad Hazan et al. *Introduction to online convex optimization*, volume 2. Now Publishers, Inc.,  
437 2016.
- 438 [29] Elad Hazan and Comandur Seshadhri. Adaptive algorithms for online decision problems. In  
439 *Electronic colloquium on computational complexity (ECCC)*, volume 14, 2007.
- 440 [30] Nicole Immorlica, Karthik Sankararaman, Robert Schapire, and Aleksandrs Slivkins. Adversar-  
441 ial bandits with knapsacks. *J. ACM*, 69(6), 2022.
- 442 [31] Nicole Immorlica, Karthik Abinav Sankararaman, Robert Schapire, and Aleksandrs Slivkins.  
443 Adversarial bandits with knapsacks. In *60th IEEE Annual Symposium on Foundations of  
444 Computer Science, FOCS 2019*, pages 202–219. IEEE Computer Society, 2019.
- 445 [32] Thomas Kesselheim and Sahil Singla. Online learning with vector costs and bandits with  
446 knapsacks. In *Conference on Learning Theory*, pages 2286–2305. PMLR, 2020.
- 447 [33] Raunak Kumar and Robert Kleinberg. Non-monotonic resource utilization in the bandits with  
448 knapsacks problem. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- 449 [34] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to  
450 personalized news article recommendation. In *Proceedings of the 19th international conference  
451 on World Wide Web*, pages 661–670, 2010.
- 452 [35] Shang Liu, Jiashuo Jiang, and Xiaocheng Li. Non-stationary bandits with knapsacks. *Advances  
453 in Neural Information Processing Systems*, 35:16522–16532, 2022.
- 454 [36] Gergely Neu. Explore no more: Improved high-probability regret bounds for non-stochastic  
455 bandits. *Advances in Neural Information Processing Systems*, 28, 2015.
- 456 [37] Francesco Orabona and Dávid Pál. Scale-free online learning. *Theoretical Computer Science*,  
457 716:50–69, 2018.

- 458 [38] David Simchi-Levi and Yunzong Xu. Bypassing the monster: A faster and simpler optimal  
459 algorithm for contextual bandits under realizability. *Mathematics of Operations Research*,  
460 47(3):1904–1931, 2022.
- 461 [39] Aleksandrs Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in*  
462 *Machine Learning*, 12(1-2):1–286, 2019.
- 463 [40] Aleksandrs Slivkins, Karthik Abinav Sankararaman, and Dylan J Foster. Contextual bandits  
464 with packing and covering constraints: A modular lagrangian approach via regression. In *The*  
465 *Thirty Sixth Annual Conference on Learning Theory*, pages 4633–4656. PMLR, 2023.
- 466 [41] Aleksandrs Slivkins, Karthik Abinav Sankararaman, and Dylan J. Foster. Contextual bandits  
467 with packing and covering constraints: A modular lagrangian approach via regression. In  
468 *The Thirty Sixth Annual Conference on Learning Theory, COLT 2023*, volume 195, pages  
469 4633–4656, 2023.
- 470 [42] Vladimir G Vovk. A game of prediction with expert advice. In *Proceedings of the eighth annual*  
471 *conference on Computational learning theory*, pages 51–60, 1995.
- 472 [43] Qian Wang, Zongjun Yang, Xiaotie Deng, and Yuqing Kong. Learning to bid in repeated  
473 first-price auctions with budgets. In *Proceedings of the 40th International Conference on*  
474 *Machine Learning, ICML’23*, 2023.
- 475 [44] Xiaohan Wei, Hao Yu, and Michael J Neely. Online primal-dual mirror descent under stochastic  
476 constraints. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*,  
477 4(2):1–36, 2020.
- 478 [45] Hao Yu, Michael Neely, and Xiaohan Wei. Online convex optimization with stochastic con-  
479 straints. *Advances in Neural Information Processing Systems*, 30, 2017.
- 480 [46] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In  
481 *Proceedings of the 20th international conference on machine learning (icml-03)*, pages 928–936,  
482 2003.

483 **A Omitted Proofs from Section 4 and Section 5**

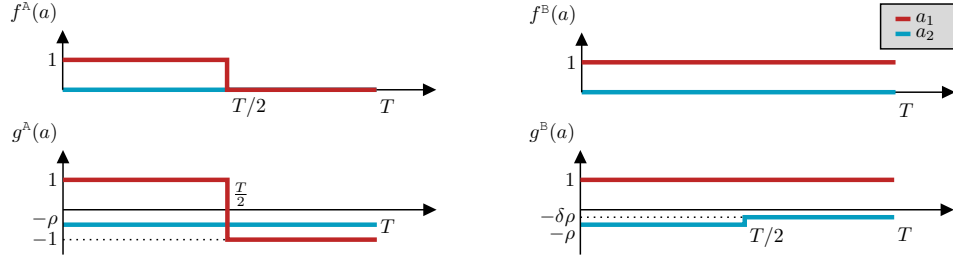


Figure 2: Lower bound adversarial setting: rewards and costs in the two instances A and B.

484 **Theorem 4.1.** [Lower bound adversarial setting] Consider the family of all adversarial instances  
 485 with  $\mathcal{X} = \{a_1, a_2\}$ , each characterized by a parameter  $\rho_{\text{Adv}}$  and optimal reward  $\text{Opt}_{\text{Adv}}$ . Then,  
 486 no algorithm can achieve, on all instances, sublinear cumulative violations  $\mathbb{E}[V(T)] = o(T)$  and  
 487  $\text{Opt}_{\text{Adv}}/\mathbb{E}[\text{Rew}] > 1 + 1/\rho_{\text{Adv}}$ .

488 *Proof.* We show that, for all  $\epsilon > 0$  and  $\delta \in (0, 1)$ , there exists two instances such that it is impossible  
 489 to obtain  $\mathbb{E}[V(T)] \leq \epsilon T$  and

$$\frac{\text{Opt}_{\text{Adv}}}{\mathbb{E}[\text{Rew}]} \geq \frac{1 + \rho_{\text{Adv}}}{\rho_{\text{Adv}}(1 + \delta) + 2\epsilon}$$

490 in both instances. The two instances are denoted by A and B respectively, with  $\mathcal{X} = \{a_1, a_2\}$  and  
 491 sequences of inputs of length  $T$ . The two instances are identical in the first  $T/2$  rounds. Rewards  
 492 in instance A are, for each  $t \in \llbracket T \rrbracket$ ,  $f_t^A(a_2) = 0$  and  $f_t^A(a_1) = \mathbb{1}[t \leq T/2]$ . On the other hand, in  
 493 instance B we have  $f_t^B(a_2) = 0$ , and  $f_t^B(a_1) = 1$  for all  $t \in \llbracket T \rrbracket$ . Costs for the first instance A are  
 494 define as

$$g_t^A(a_1) := \begin{cases} 1 & \text{if } t \leq T/2 \\ -1 & \text{otherwise} \end{cases},$$

495 and  $g_t^A(a_2) = -\rho$  for all  $t \in \llbracket T \rrbracket$ . In the second instance B, costs are  $g_t^B(a_1) = 1$  for all  $t \in \llbracket T \rrbracket$ , and

$$g_t^B(a_2) := \begin{cases} -\rho & \text{if } t \leq T/2 \\ -\delta\rho & \text{otherwise} \end{cases},$$

496 for some  $\delta > 0$ . The two instances are depicted in Figure 2.

497 Let  $N$  be the expected number of times that action  $a_1$  is played in rounds  $\llbracket T/2 \rrbracket$ , that is

$$N := \sum_{t \in \llbracket T/2 \rrbracket} \mathbb{E}^A[x_t = a_1] = \sum_{t \in \llbracket T/2 \rrbracket} \mathbb{E}^B[x_t = a_1],$$

498 where expectation is with respect to the algorithm's randomization. We observe that the algorithm  
 499 plays in the same way in both instances up to time  $T/2$ , as they are identical (formally, the KL  
 500 between instance A and B is zero in the first  $T/2$  rounds). Then, we have that the optimal action in  
 501 instance A is to play deterministically action  $a_1$ . Therefore,  $\text{Opt}_{\text{Adv}}^A = T/2$ . The expected reward in  
 502 instance A comes only from the number of plays of  $a_1$  in the first  $T/2$  rounds:  $\mathbb{E}^A[\text{Rew}] = N$ . On the  
 503 other hand, call  $M$  the expected number of times an algorithm plays action  $a_1$  in the last  $\llbracket T/2, T \rrbracket$   
 504 rounds of instance B, that is

$$M := \sum_{t \in \llbracket T/2, T \rrbracket} \mathbb{E}^B[x_t = a_1].$$

505 We have that, in order to have  $\mathbb{E}^B[V(T)] \leq \epsilon T$  violations in the second instance, we need to play  $a_1$   
 506 a small number of times:

$$M - \delta\rho\left(\frac{T}{2} - M\right) + N - \rho\left(\frac{T}{2} - N\right) \leq \epsilon T,$$

507 which yields

$$N \leq \frac{T(\rho(\delta + 1) + 2\epsilon)}{2(\rho + 1)}.$$

508 Then, we get that

$$\frac{\text{Opt}_{\text{Adv}}^{\text{A}}}{\mathbb{E}^{\text{A}}[\text{Rew}]} \geq \frac{1 + \rho}{\rho(1 + \delta) + 2\epsilon},$$

509 which concludes the proof since  $\rho_{\text{Adv}}^{\text{A}} = \rho$ .  $\square$

510 **Proposition 5.3.** *Consider the instance of Example 5.2. Even if  $\text{A1}\mathcal{G}_{\text{P}}$  and  $\text{A1}\mathcal{G}_{\text{D}}$  suffer regret less*  
 511 *than or equal then zero, the primal-dual framework fails to achieve sublinear constraint violations.*

512 *Proof.* Consider the instance described in Example 5.2, and consider an algorithm  $\text{A1}\mathcal{G}_{\text{P}}$  for  $\mathcal{X} =$   
 513  $\{a_1, a_2, a_3\}$  such that  $x_t = a_3$  for  $t \in \llbracket T/3 \rrbracket$ , while  $x_t = a_2$  for  $t \in \llbracket T/3, T \rrbracket$ . Moreover, consider  
 514 an algorithm  $\text{A1}\mathcal{G}_{\text{D}}$  instantiated on  $\mathcal{D} = [0, M]$ , with  $M \geq 1/\rho$ , that plays  $\lambda_t = 0$  for all  $t \in \llbracket 2T/3 \rrbracket$ ,  
 515 and  $\lambda_t = M$  for all  $t \in \llbracket 2T/3, T \rrbracket$ .

516 We start by analyzing the primal regret achieved by  $\text{A1}\mathcal{G}_{\text{P}}$ :

$$\begin{aligned} R_T^{\text{P}} &:= \sup_{x \in \mathcal{X}} \sum_{t \in \llbracket T \rrbracket} [f_t(x) - f_t(x_t) - \lambda_t(g_{t,1}(x) - g_{t,1}(x_t))] \\ &= \sup_{x \in \mathcal{X}} \sum_{t \in \llbracket T \rrbracket} [f_t(x) - \lambda_t g_{t,1}(x)] - \frac{2}{3}T + \frac{M\rho}{3}T \\ &= \sum_{t \in \llbracket T \rrbracket} [f_t(a_1) - \lambda_t g_{t,1}(a_1)] + \frac{T}{3}(M\rho - 2) \\ &= \rho M \frac{T}{3} + \frac{T}{3}(M\rho - 2) \\ &= \frac{T}{3}(2M\rho - 2) \leq 0, \end{aligned}$$

517 where we replaced the sup with the utility at  $a_1$  since  $M \geq 1/\rho$ . Moreover, the dual regret is such  
 518 that

$$\begin{aligned} R_T^{\text{D}} &:= \sup_{\lambda \in [0, M]} \sum_{t \in \llbracket 2T/3, T \rrbracket} (\lambda - M) g_{t,1}(x_t) \\ &= \sup_{\lambda \in [0, M]} \frac{T}{3} (\lambda - M) \rho = 0. \end{aligned}$$

519 However, for a suitable choice of  $\rho$ , the violations are linear in  $T$  since

$$V_1(T) := \sum_{t \in \llbracket T \rrbracket} g_{t,1}(x_t) = \frac{\rho}{3}T = \Omega(T).$$

520 This concludes the proof.  $\square$

## 521 B Proof of Lemma 6.2

522 We start by providing the following auxiliary lemmas.

523 **Lemma B.1.** *Let  $\mathbf{y}_t \in \mathbb{R}_{\geq 0}^m$  be generated by OGD with learning rate  $\eta$  and utilities  $\mathbf{y} \mapsto \langle \mathbf{y}, \mathbf{g}_t \rangle$ ,*  
 524 *where  $\|\mathbf{g}_t\|_{\infty} \leq 1$  for all  $t \in \llbracket T \rrbracket$ . Then:*

$$\|\mathbf{y}_{t+1}\|_1 - \|\mathbf{y}_t\|_1 \leq m \cdot \eta$$

525 *Proof.* The update of the  $i$ -th component of  $\mathbf{y}_{t+1}$  can be written as:

$$y_{t+1,i} := \max(0, y_{t,i} + \eta g_{t,i}).$$

526 If  $g_{t,i} \geq 0$  then the update can be simplified to  $y_{t+1,i} = y_{t,i} + \eta g_{t,i} \leq y_{t,i} + \eta$ . If  $g_{t,i} < 0$  then  
 527  $y_{t+1,i} \geq y_{t,i} + \eta g_{t,i} \geq y_{t,i} - \eta$ . Thus  $|y_{t+1,i} - y_{t,i}| \leq \eta$  for all  $i \in \llbracket m \rrbracket$ . By summing over  
 528 all component we have that  $\|\mathbf{y}_{t+1} - \mathbf{y}_t\|_1 \leq m \cdot \eta$ . By triangular inequality we have the desired  
 529 statement.  $\square$

530 **Lemma B.2.** *[[28, Chapter 10]] For any  $t_1, t_2 \in \llbracket T \rrbracket$  with  $t_1 < t_2$ , it holds that if  $\lambda_t$  is generated*  
 531 *by OGD with learning rate  $\eta > 0$  on a set  $\mathcal{D}$ , then:*

$$R_{\llbracket t_1, t_2 \rrbracket}^p(\{\lambda\}) \leq \frac{\|\lambda - \lambda_{t_1}\|_2^2}{2\eta} + \frac{1}{2}\eta mT.$$

532 *with probability one on the randomization of the algorithm, i.e.,  $\delta = 0$ . Moreover it also*  
 533 *holds component-wise, i.e., for all  $\lambda \geq 0$ :*

$$\sum_{t \in \llbracket t_1, t_2 \rrbracket} (\lambda - \lambda_t) g_t(x_t) \leq \frac{(\lambda - \lambda_{t_1})^2}{2\eta} + \frac{1}{2}\eta T.$$

534 **Lemma B.3.** *In the stochastic setting, for any  $\xi \in \Delta(\mathcal{X})$  and  $\delta \in (0, 1]$ , with probability at least*  
 535  *$1 - \delta$ , it holds that:*

$$\sum_{t \in I} \mathbb{E}_{x \sim \xi} [\langle \lambda_t, \mathbf{g}_t(x) \rangle] \leq \sum_{t \in I} \mathbb{E}_{x \sim \xi} [\langle \lambda_t, \bar{\mathbf{g}}_t(x) \rangle] + ME_{T,\delta} \quad \text{and} \quad (1)$$

$$\sum_{t \in I} \mathbb{E}_{x \sim \xi} [f_t(x)] \geq \sum_{t \in I} \mathbb{E}_{x \sim \xi} [\bar{f}(x)] - E_{T,\delta}, \quad (2)$$

536 *for any interval  $I = [t_1, t_2] \subseteq [T]$ , where  $E_{T,\delta} := \sqrt{16T \log\left(\frac{2T}{\delta}\right)}$  and  $M = \sup_{t \in \llbracket T \rrbracket} \|\lambda\|_1$ .*

537 *Proof.* We start by proving that the all the inequalities of Equation (1) holds simultaneously with  
 538 *probability  $1 - \delta/2$ . We have that given a  $I = [t_1, t_2] \subseteq [T]$ , with probability at least  $1 - \delta/(2T^2)$ ,*

$$\sum_{t \in I} \mathbb{E}_{x \sim \xi} [\langle \lambda_t, \mathbf{g}_t(x) \rangle] - \sum_{t \in I} \mathbb{E}_{x \sim \xi} [\langle \lambda_t, \bar{\mathbf{g}}_t(x) \rangle] \leq M \sqrt{8|I| \log\left(\frac{2T^2}{\delta}\right)} \leq M \sqrt{16T \log\left(\frac{2T}{\delta}\right)},$$

539 *where the first inequality holds by Azuma-Hoeffding inequality. By taking a union bound over all*  
 540 *possible intervals  $I$  (which are at most  $T^2$ ), we obtain that all the first set of equations holdswith*  
 541 *probability at least  $1 - \delta/2$ .*

542 *Equation (2) can be proved in a similar way. Indeed, for any fixed interval  $I = [t_1, t_2] \subseteq [T]$ , and for*  
 543 *any strategy mixture  $\xi \in \Delta(\mathcal{X})$ , by the Azuma-Hoeffding inequality we have that, with probability*  
 544 *at least  $1 - \delta/(2T^2)$ , the following holds*

$$\sum_{t \in I} \mathbb{E}_{x \sim \xi} [\bar{f}(x)] - \sum_{t \in I} \mathbb{E}_{x \sim \xi} [f_t(x)] \leq \sqrt{2|I| \log\left(\frac{2T^2}{\delta}\right)} \leq \sqrt{4T \log\left(\frac{2T}{\delta}\right)}.$$

545 *By taking a union bound over all possible  $T^2$  intervals, we obtain that, for all possible intervals  $I$ , the*  
 546 *equation above holds with probability  $1 - \delta/2$ .*

547 *The Lemma follows by a union bound on the two sets of equations above.  $\square$*

548 *These auxiliary technical lemmas are used in proving the following result.*

549 **Lemma 6.2 (Self-bounding lemma).** *Let  $\eta_{\text{OGD}} := (800 \cdot m \cdot \max\{\bar{R}_{T,\delta}^p(\mathcal{X}), E_{T,\delta}\})^{-1}$ , then if*  
 550  *$\text{Alg}_{\mathcal{D}}$  is OGD on the set  $\mathcal{D} = \mathbb{R}_{\geq 0}^m$ , and the primal algorithm  $\text{Alg}_{\mathcal{P}}$  is 2-scale-free and has a*  
 551 *high-probability weakly adaptive regret bound  $\bar{R}_{T,\delta}^p(\mathcal{X})$ , then with probability at least  $1 - \delta$ :*

$$\max_{t \in \llbracket T \rrbracket} \|\lambda_t\|_1 \leq \frac{13m}{\rho},$$

552 *where  $\rho = \rho_{\text{Adv}}$  or  $\rho = \rho_{\text{Stoc}}$  depending on the setting and  $E_{T,\delta} := \sqrt{16T \log(2T/\delta)}$ .*

553 *Proof.* Let  $c_1 := 2$  and  $c_2 := 12m$  and any learning rate  $\eta$  for OGD with  $\eta \leq \eta_{\text{OGD}}$ . By contradiction,  
 554 *suppose there exists a time such that  $\|\lambda_t\|_1 \geq c_2/\rho$ , and let  $t_2 \in \llbracket T \rrbracket$  be the smallest  $t$  for which this*  
 555 *happens. We unify the proof of the adversarial and stochastic setting. In particular, let  $\rho = \rho_{\text{Adv}}$  if*  
 556 *the losses  $(f_t, \mathbf{g}_t)$  are adversarial, and let  $\rho = \rho_{\text{Stoc}}$  if  $(f_t, \mathbf{g}_t)$  are stochastic with mean  $(\bar{f}, \bar{\mathbf{g}})$ . The*  
 557 *extra stochasticity coming from the environment in the stochastic setting will be handled through*  
 558 *Lemma B.3. In order to streamline the notation, we define  $E_{T,\delta} := \sqrt{16T \log(2T/\delta)}$ .*

559 Then, let  $t_1 \in \llbracket t_2 \rrbracket$  be the largest time for which  $\|\lambda_t\|_1 \in [\frac{c_1}{\rho}, \frac{c_2}{\rho}]$  for all  $t \in \llbracket t_1, t_2 \rrbracket$ .

560 **Step 1.** First, we need to bound  $\|\lambda_{t_1}\|_1$  and  $\|\lambda_{t_2}\|_1$ . To do that, we exploit Lemma B.1. In particular,  
561 by telescoping the sum in the lemma, we obtain that:

$$\|\lambda_{t_2}\|_1 - \|\lambda_{t_1}\|_1 \leq \eta m(t_2 - t_1).$$

562 Moreover, by the definition of  $\lambda_{t_1}$  and  $\lambda_{t_2}$ , we have:

$$\frac{c_1}{\rho} \leq \|\lambda_{t_1}\|_1 \leq \|\lambda_{t_1-1}\|_1 + m\eta \leq \frac{c_1}{\rho} + m\eta$$

563 and similarly

$$\frac{c_2}{\rho} \leq \|\lambda_{t_2}\|_1 \leq \|\lambda_{t_2-1}\|_1 + m\eta \leq \frac{c_2}{\rho} + m\eta.$$

564 This, together with the inequality above, yields

$$\frac{c_2 - c_1}{2\eta m \rho} \leq t_2 - t_1. \quad (3)$$

565 **Step 2.** The range of the primal utilities in the turns  $\llbracket t_1, t_2 \rrbracket$  can now be bounded as:

$$\begin{aligned} \sup_{x \in \mathcal{X}, t \in \llbracket t_1, t_2 \rrbracket} |u_t^p(x)| &\leq \sup_{x \in \mathcal{X}, t \in \llbracket t_1, t_2 \rrbracket} \{|f_t(x)| + \|\lambda_t\|_1 \cdot \|\mathbf{g}_t(x)\|_\infty\} \\ &\leq 1 + \frac{c_2}{\rho} + m\eta \\ &\leq 1 + \frac{12m + 1}{\rho} \\ &\leq \frac{14m}{\rho} =: L. \end{aligned}$$

566 Now, by the assumption that  $\text{Alg}_p$  is weakly adaptive and 2-scale-free, we obtain:

$$R_{\llbracket t_1, t_2 \rrbracket}^p(\mathcal{X}) \leq L^2 \cdot \overline{R}_{T, \delta}^p(\mathcal{X}),$$

567 which holds with probability at least  $1 - \delta$ .

568 If we apply the primal no-regret condition above for strictly safe strategy  $\xi^\circ \in \Delta(\mathcal{X})$  we have

$$\sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathcal{L}_{f_t, \mathbf{g}_t}(x_t, \lambda_t) \geq \mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathcal{L}_{f_t, \mathbf{g}_t}(x, \lambda_t) \right] - L^2 \overline{R}_{T, \delta}^p(\mathcal{X}). \quad (4)$$

569 Moreover, by definition of safe strategy we have that in the adversarial setting  $\mathbb{E}_{x \sim \xi^\circ} [g_{t,i}(x)] \leq$   
570  $-\rho_{\text{Adv}}$  for all  $i \in \llbracket m \rrbracket$  and  $t \in \llbracket t_1, t_2 \rrbracket$ , while in the stochastic setting by Lemma B.3 it holds

$$\sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathbb{E}_{x \sim \xi^\circ} [\langle \lambda_t, \mathbf{g}_t(x) \rangle] \leq \sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathbb{E}_{x \sim \xi^\circ} [\langle \lambda_t, \bar{\mathbf{g}}_t(x) \rangle] + M E_{T, \delta}$$

571 and

$$\mathbb{E}_{x \sim \xi^\circ} [\bar{g}_i(\xi)] \leq -\rho_{\text{Stoc}} \quad \forall i \in \llbracket m \rrbracket,$$

572 where we recall that  $E_{T, \delta} = \sqrt{16T \log(2T/\delta)}$  and  $M = \sup_{t \in \llbracket T \rrbracket} \|\lambda\|_1$ .



573 Therefore, we can lower bound the first term of the right-hand side of Equation (4) the stochastic  
 574 setting as:

$$\begin{aligned}
 \mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathcal{L}_{f_t, g_t}(x, \lambda_t) \right] &= \mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in \llbracket t_1, t_2 \rrbracket} f_t(x) - \langle \lambda_t, g_t(x) \rangle \right] \\
 &\geq -\mathbb{E}_{x \sim \xi^\circ} [\langle \lambda_t, g_t(x) \rangle] \\
 &\geq -\mathbb{E}_{x \sim \xi^\circ} [\langle \lambda_t, \bar{g}(x) \rangle] - \left( \sup_{t \in \llbracket T \rrbracket} \|\lambda\|_1 \right) E_{T, \delta} \\
 &\geq \rho_{\text{Stoc}} \sum_{t \in \llbracket t_1, t_2 \rrbracket} \|\lambda_t\|_1 - \left( \sup_{t \in \llbracket T \rrbracket} \|\lambda\|_1 \right) E_{T, \delta} \\
 &\geq \rho_{\text{Stoc}} \sum_{t \in \llbracket t_1, t_2 \rrbracket} \|\lambda_t\|_1 - \left( \frac{c_2}{\rho_{\text{Stoc}}} + m\eta \right) E_{T, \delta} \\
 &\geq c_1(t_2 - t_1) - \left( \frac{c_2}{\rho_{\text{Stoc}}} + m\eta \right) E_{T, \delta}
 \end{aligned}$$

575 In the adversarial setting we can more easily conclude that  $\mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathcal{L}_{f_t, g_t}(x, \lambda_t) \right] \geq c_1(t_2 -$   
 576  $t_1)$  and thus in both settings it holds that:

$$\mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathcal{L}_{f_t, g_t}(x, \lambda_t) \right] \geq c_1(t_2 - t_1) - \left( \frac{c_2}{\rho_{\text{Stoc}}} + m\eta \right) E_{T, \delta}. \quad (5)$$

577 Combining the two inequalities of Equation (4) and Equation (5), we can conclude that the overall  
 578 utility of the primal algorithm  $\text{Al}_{\mathcal{G}_P}$  can be lower bounded by:

$$\sum_{t \in \llbracket t_1, t_2 \rrbracket} u_t^P(x_t) \geq c_1(t_2 - t_1) - L^2 \overline{R}_{T, \delta}^P(\mathcal{X}) - \left( \frac{c_2}{\rho} + m\eta \right) E_{T, \delta} \quad (6)$$

579 Now, we need an auxiliary result that we will use to upper bound the left hand side of the previous  
 580 inequality.

581 **Claim B.4.** *It holds that:*

$$\sum_{t \in \llbracket t_1, t_2 \rrbracket} \langle \lambda_t, g_t(x_t) \rangle \geq \frac{m}{2\rho^2\eta}.$$

582 Then, we upper bound the left-hand side by using Claim B.4:

$$\begin{aligned}
 \sum_{t \in \llbracket t_1, t_2 \rrbracket} u_t^P(x_t) &= \sum_{t \in \llbracket t_1, t_2 \rrbracket} \mathcal{L}_{f_t, g_t}(x_t, \lambda_t) = \sum_{t \in \llbracket t_1, t_2 \rrbracket} [f_t(x_t) - \langle \lambda_t, g_t(x_t) \rangle] \\
 &\leq (t_2 - t_1) - \frac{m}{2\rho^2\eta}
 \end{aligned} \quad (7)$$

583 Thus, combining Equation (7) and (6)

$$t_2 - t_1 \leq \frac{1}{c_1 - 1} \left( L^2 \overline{R}_{T, \delta}^P(\mathcal{X}) - \frac{m}{2\rho^2\eta} + \left( \frac{c_2}{\rho} + m\eta \right) E_{T, \delta} \right).$$

584 Combining it with Equation (3) one obtains that:

$$\frac{c_2 - c_1}{2\eta m \rho} \leq \frac{1}{c_1 - 1} \left( L^2 \overline{R}_{T, \delta}^P(\mathcal{X}) - \frac{m}{2\rho^2\eta} + \left( \frac{c_2}{\rho} + m\eta \right) E_{T, \delta} \right),$$

585 which gives as a solution  $\eta \geq \frac{m^2 - 2\rho + 13m\rho}{392m^3 \bar{R}_{T,\delta}^2 (\mathcal{X} + 2m\rho E_{T,\delta} (1 + 13m))}$ . Which is a contradiction since:

$$\eta \leq \eta_{\text{GD}} := \frac{1}{800 \cdot m \cdot \max\{\bar{R}_{T,\delta}^2(\mathcal{X}), E_{T,\delta}\}} > \frac{m^2 - 2\rho + 13m\rho}{392m^3 \bar{R}_{T,\delta}^2 (\mathcal{X} + 2m\rho E_{T,\delta} (1 + 13m))}$$

586 Thus, we can conclude that  $\|\lambda_t\|_t \leq c_2/\rho$  for each  $t \in \llbracket T \rrbracket$ .  $\square$

587 Now, we provide the proof of Claim B.4.

588 **Proof of Claim B.4.** We define  $\tilde{t}_i$  as the last time in  $\llbracket t_1, t_2 \rrbracket$  in which  $\lambda_{\tilde{t}_i, i} = 0$ , or  $\tilde{t}_i = t_1$  if  
589  $\lambda_{t_1, i} > 0$  for all  $t \in \llbracket t_1, t_2 \rrbracket$ . Formally:

$$\tilde{t}_{1,i} = \max \left\{ t_1, \sup_{\tau \in \llbracket t_2 \rrbracket : \lambda_{\tau, i} = 0} \tau \right\}.$$

590 We are now going to analyze separately for all  $i \in \llbracket m \rrbracket$ , the rounds  $\llbracket t_1, \tilde{t}_{1,i} \rrbracket$  and the rounds  $\llbracket \tilde{t}_{1,i}, t_2 \rrbracket$ .

591 **Phase 1:** First, we analyze the rounds  $\llbracket t_1, \tilde{t}_{1,i} \rrbracket$ . By definition, it can be either that  $\lambda_{\tilde{t}_{1,i}, i} = 0$  or  
592  $\tilde{t}_{1,i} = t_1$ . In the latter case,  $\llbracket t_1, \tilde{t}_{1,i} \rrbracket = \emptyset$  and the dual algorithm incurs zero regret. In the former  
593 case, we can use Lemma B.2 and write that the regret over the interval with respect to  $\lambda_i^* = 0$  is

$$0 \leq \sum_{t \in \llbracket t_1, \tilde{t}_{1,i} \rrbracket} \lambda_{t,i} g_{t,i}(x_t) + \frac{\lambda_{t_1}^2}{2\eta} + \frac{1}{2}\eta T \leq \sum_{t \in \llbracket t_1, \tilde{t}_{1,i} \rrbracket} \lambda_{t,i} g_{t,i}(x_t) + \frac{\lambda_{t_1}^2}{2\eta} + \frac{1}{2}\eta T. \quad (8)$$

594 **Phase 2:** Now, we consider the rounds  $\llbracket \tilde{t}_{1,i}, t_2 \rrbracket$ . We take  $\lambda^*$  defined as follows:  $\lambda_i^* = \frac{1}{\rho}$  for all  
595  $i \in \llbracket m \rrbracket$ .

596 Let  $\tilde{\Delta}_i := \lambda_{t_2, i} - \lambda_{\tilde{t}_{1,i}, i}$ . Due to the definition of  $\tilde{t}_{1,i}$ , gradient descent never projects the multiplier  
597 relative to constraint  $i$ , and we can write that

$$\sum_{t \in \llbracket \tilde{t}_{1,i}, t_2 \rrbracket} g_{t,i}(x_t) = \frac{\tilde{\Delta}_i}{\eta}$$

598 and, therefore,

$$\sum_{t \in \llbracket \tilde{t}_{1,i}, t_2 \rrbracket} \lambda_i^* g_{t,i}(x_t) = \frac{\tilde{\Delta}_i}{\rho\eta}. \quad (9)$$

599 Now we can use Lemma B.2 to find that:

$$\sum_{t \in \llbracket \tilde{t}_{1,i}, t_2 \rrbracket} \lambda_i^* g_{t,i}(x_t) \leq \sum_{t \in \llbracket \tilde{t}_{1,i}, t_2 \rrbracket} \lambda_{t,i} g_{t,i}(x_t) + \frac{(\lambda_i^* - \lambda_{\tilde{t}_{1,i}, i})^2}{2\eta} + \frac{1}{2}\eta T.$$

600 Combining it with Equation (9) yields the following

$$\sum_{t \in \llbracket \tilde{t}_{1,i}, t_2 \rrbracket} \lambda_{t,i} g_{t,i}(x_t) \geq \frac{\tilde{\Delta}_i}{\rho\eta} - \frac{(\lambda_i^* - \lambda_{\tilde{t}_{1,i}, i})^2}{2\eta} - \frac{1}{2}\eta T. \quad (10)$$

601 Combining Equation (10) and Equation (8) we obtain:

$$\begin{aligned} \sum_{t \in \llbracket t_1, t_2 \rrbracket} \lambda_{t,i} g_{t,i}(x_t) &\geq \frac{\tilde{\Delta}_i}{\rho\eta} - \frac{(\lambda_i^* - \lambda_{\tilde{t}_{1,i}, i})^2}{2\eta} - \frac{\lambda_{t_1}^2}{2\eta} - \eta T \\ &\geq \frac{\tilde{\Delta}_i}{\rho\eta} - \frac{(\lambda_i^*)^2 + \lambda_{\tilde{t}_{1,i}, i}^2}{2\eta} - \frac{\lambda_{t_1}^2}{2\eta} - \eta T. \end{aligned}$$

602 Now, by summing over all  $i \in \llbracket m \rrbracket$ , and by letting  $\lambda_{\tilde{t}_1}$  be the vector that has  $\lambda_{\tilde{t}_1, i}$  as its  $i$ -th  
 603 component, we get:

$$\begin{aligned}
 \sum_{t \in \llbracket t_1, t_2 \rrbracket} \langle \lambda_t, g_t(x_t) \rangle &\geq \frac{\|\lambda_{t_2}\|_1 - \|\lambda_{\tilde{t}_1}\|_1}{\rho\eta} - \frac{1}{2\eta} (\|\lambda^*\|_2^2 + \|\lambda_{\tilde{t}_1}\|_2^2 + \|\lambda_{t_1}\|_2^2) - \frac{1}{\eta} \quad (\text{as } \eta \leq 1/\sqrt{T}) \\
 &\geq \frac{c_2}{\rho^2\eta} - \frac{1}{\rho\eta} \|\lambda_{t_1}\|_1 - \frac{1}{2\eta} (\|\lambda^*\|_2^2 + 2\|\lambda_{t_1}\|_2^2) - \frac{1}{\eta} \\
 &\hspace{15em} (\|\lambda\|_1 \geq c_2/\rho \text{ and } \|\lambda_{\tilde{t}_1}\|_1 \leq \|\lambda_{t_1}\|_1) \\
 &\geq \frac{c_2}{\rho^2\eta} - \frac{1}{\rho\eta} \left( \frac{c_1}{\rho} + m\eta \right) - \frac{1}{2\eta} \left( \frac{m}{\rho^2} + 2 \left( \frac{c_1}{\rho} + m\eta \right)^2 \right) - \frac{1}{\eta} \\
 &\geq \frac{c_2}{\rho^2\eta} - \frac{c_1 + 1}{\rho^2\eta} - \frac{m}{2\rho^2\eta} - \frac{2(c_1 + 1)^2}{2\rho^2\eta} - \frac{1}{\eta} \quad (\eta \leq 1/\rho m) \\
 &\geq \frac{2c_2 - 24 - m}{2\rho^2\eta} \\
 &\geq \frac{m}{2\rho^2\eta}
 \end{aligned}$$

604 where the last two inequalities hold due to the choice of parameters in the proof of Claim B.4, that is  
 605  $c_1 = 2$  and  $c_2 = 13m$ . This concludes the proof.  $\square$

## 606 C Omitted Proofs from Section 7

607 **Theorem 7.1.** *Let  $\text{Alg}_{\mathcal{D}}$  be OGD with learning rate  $\eta$  as in Lemma 6.2, and let  $\text{Alg}_{\mathcal{P}}$  any 2-  
 608 scale-free algorithm with no-adaptive regret. Then, with probability at least  $1 - \delta$ , it holds that  
 609  $V_T = \tilde{O} \left( \frac{m^2}{\rho} \overline{R}_{T, \delta}^{\mathcal{P}}(\mathcal{X}) \right)$ , where  $\rho = \rho_{\text{Adv}}$  in the adversarial setting and  $\rho = \rho_{\text{Stoc}}$  in the stochastic.*

610 *Proof.* The update of OGD for each component  $i \in \llbracket m \rrbracket$  is  $\lambda_{t+1, i} := [\lambda_{t, i} + \eta g_{t, i}(x_t)]^+$ . Thus:

$$\lambda_{t+1, i} \geq \lambda_{t, i} + \eta_{\text{OGD}} g_{t, i}(x_t),$$

611 and by induction:

$$\lambda_{t+1, i} \geq \lambda_{0, i} + \eta_{\text{OGD}} \sum_{\tau=1}^t g_{\tau, i}(x_{\tau}).$$

612 By rearranging and recalling that  $\lambda_{0, i} = 0$  we obtain:

$$\sum_{t \in \llbracket T \rrbracket} g_{t, i}(x_t) \leq \frac{1}{\eta_{\text{OGD}}} \lambda_{T+1, i} \leq \frac{1}{\eta} \|\lambda_{T+1}\|_1$$

613 Moreover, by Lemma 6.2 we can bound  $\|\lambda_T\|_1 \leq \frac{13m}{\rho}$  which holds with probability at least  $1 - \delta$ .

614 Thus, with probability at least  $1 - \delta$ , it holds:

$$V_T := \max_{i \in \llbracket m \rrbracket} V_i(T) \leq \frac{13m}{\eta_{\text{OGD}} \rho}.$$

615 The proof is concluded by observing that  $\eta_{\text{OGD}} = \tilde{O} \left( (m \overline{R}_{T, \delta}^{\mathcal{P}}(\mathcal{X}))^{-1} \right)$ .  $\square$

616 **Theorem 7.2.** *If  $\text{Alg}_{\mathcal{D}}$  is OGD with learning rate  $\eta_{\text{OGD}}$  and domain  $\mathcal{D} := \mathbb{R}_{\geq 0}^m$ , and  $\text{Alg}_{\mathcal{P}}$  is 2-scale-  
 617 free, then, in the adversarial setting, with high probability:*

$$\text{Re}w \geq \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} - \tilde{O} \left( \left( \frac{m}{\rho_{\text{Adv}}} \right)^2 \overline{R}_{T, \delta}^{\mathcal{P}}(\mathcal{X}) \right).$$

618 *Proof.* Define  $x^* \in \mathcal{X}$  such that:

$$\sum_{t \in [T]} f_t(x^*) = \text{Opt}_{\text{Adv}}$$

619 Now, consider a randomized strategy  $\xi$  that randomized with probability  $\alpha$  between  $x^*$  and  $\xi^\circ$ , where  
 620  $\xi^\circ$  is any strategy for which  $\mathbb{E}_{x \sim \xi^\circ} [g_{t,i}(x_t)] \leq -\rho_{\text{Adv}}$ . This strategy exists by assumption. Formally,  
 621 for any  $x \in \mathcal{X}$  the randomized strategy  $\xi$  assigns probability to  $x$ :

$$\xi(x) = \alpha \delta_{x^*}(x) + (1 - \alpha) \xi^\circ(x).$$

622 Then, we compute the component of the primal utility of  $\xi$  due to a constraint  $i \in [m]$  as follows:

$$\begin{aligned} \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} \lambda_{t,i} g_{t,i}(x) \right] &= \alpha \sum_{t \in [T]} \lambda_{t,i} g_{t,i}(x^*) + (1 - \alpha) \mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in [T]} \lambda_{t,i} g_{t,i}(x) \right] \\ &\leq \alpha \sum_{t \in [T]} \lambda_{t,i} - (1 - \alpha) \rho_{\text{Adv}} \sum_{t \in [T]} \lambda_{t,i} \\ &\leq (\alpha - (1 - \alpha) \rho_{\text{Adv}}) \sum_{t \in [T]} \lambda_{t,i}. \end{aligned}$$

623 Thus, setting  $\alpha = \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}}$  we have that  $\mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} \lambda_{t,i} g_{t,i}(x) \right] \leq 0$ , and  $\sum_{t \in [T]} \langle \lambda_t, \mathbf{g}_t(x_t) \rangle \leq 0$ .

624 We now compute the reward of  $\xi$  for  $\alpha = \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}}$ :

$$\begin{aligned} \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} f_t(x) \right] &= \alpha \sum_{t \in [T]} f_t(x^*) + (1 - \alpha) \mathbb{E}_{x \sim \xi^\circ} \left[ \sum_{t \in [T]} f_t(x) \right] \\ &\geq \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} \end{aligned}$$

625 Now, we consider the regret of  $\text{Alg}_{\text{P}}$  with respect to  $\xi$  and we find that:

$$\sum_{t \in [T]} \mathcal{L}_{f_t, \mathbf{g}_t}(x_t, \lambda_t) \geq \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} \mathcal{L}_{f_t, \mathbf{g}_t}(x, \lambda_t) \right] - L^2 \cdot \overline{R}_{T, \delta}^{\text{P}}(\mathcal{X}).$$

626 where  $L$  is the maximum module of the payoffs of the primal regret minimizer, *i.e.*,  $L :=$   
 627  $\sup_{t \in [T], x \in \mathcal{X}} |u_t^{\text{P}}(x)|$ .

628 Exploiting the definition of  $\mathcal{L}_{f_t, \mathbf{g}_t}(\cdot, \cdot)$  in the inequality above we obtain that:

$$\begin{aligned} \sum_{t \in [T]} f_t(x_t) - \langle \lambda_t, \mathbf{g}_t(x_t) \rangle &\geq \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} f_t(x) - \langle \lambda_t, \mathbf{g}_t(x) \rangle \right] - L^2 \cdot \overline{R}_{T, \delta}^{\text{P}}(\mathcal{X}) \\ &\geq \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} f_t(x) \right] - L^2 \cdot \overline{R}_{T, \delta}^{\text{P}}(\mathcal{X}) \\ &\geq \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} - L^2 \cdot \overline{R}_{T, \delta}^{\text{P}}(\mathcal{X}) \end{aligned} \quad (11)$$

629 Then, we lower bound the term  $\sum_{t \in [T]} \langle \lambda_t, \mathbf{g}_t(x_t) \rangle$  by using the dual regret of  $\text{Alg}_{\text{D}}$  with respect to

630  $\lambda^* = \mathbf{0}$ . Indeed,

$$\sum_{t \in [T]} \langle \lambda^* - \lambda_t, \mathbf{g}_t(x_t) \rangle \leq \overline{R}_{T, \delta}^{\text{D}}(\{\lambda^*\})$$

631 implies that

$$\sum_{t \in [T]} \langle \lambda_t, \mathbf{g}_t(x_t) \rangle \geq -\overline{R}_{T, \delta}^{\text{D}}(\{\lambda^*\}).$$

632 Combining it with Equation (11) gives:

$$\sum_{t \in [T]} f_t(x_t) \geq \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} - L^2 \cdot \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) - \overline{R}^{\mathbb{D}}_{T,\delta}(\{\boldsymbol{\lambda}^*\}).$$

633 Now, we use Lemma 6.2 which bounds  $L \leq 2 \frac{13m}{\rho_{\text{Adv}}}$  and Lemma B.1 which we can use to bound

634  $\overline{R}^{\mathbb{D}}_{T,\delta}(\{\boldsymbol{\lambda}^*\})$ .

635 In particular,  $\overline{R}^{\mathbb{D}}_{T,\delta}(\{\boldsymbol{\lambda}^*\})$  can be bounded with:

$$\overline{R}^{\mathbb{D}}_{T,\delta}(\{\boldsymbol{\lambda}^*\}) \leq \frac{1}{2} \eta_{\text{OGD}} mT,$$

636 and thus:

$$\text{Rew} := \sum_{t \in [T]} f_t(x_t) \geq \frac{\rho_{\text{Adv}}}{1 + \rho_{\text{Adv}}} \text{Opt}_{\text{Adv}} - 676 \left( \frac{m}{\rho_{\text{Adv}}} \right)^2 \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) - \eta_{\text{OGD}} mT.$$

637 The proof is concluded by noting that  $\eta_{\text{OGD}} = \tilde{O}((m \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}))^{-1})$ .  $\square$

638 **Theorem 7.3.** *If  $\text{Alg}_{\mathbb{D}}$  is OGD with learning rate  $\eta_{\text{OGD}}$  and domain  $\mathcal{D} := \mathbb{R}_{\geq 0}^m$ , and  $\text{Alg}_{\mathbb{P}}$  is 2-scale-*  
 639 *free, then in the stochastic setting, in high probability:*

$$\text{Rew} \geq \text{Opt}_{\text{Stoc}} - \tilde{O} \left( \left( \frac{m}{\rho_{\text{Stoc}}} \right)^2 \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) \right).$$

640 *Proof.* By Lemma 6.2 we have that with probability at least  $1 - \delta$  we have that  $\sup_{t \in [T]} \|\boldsymbol{\lambda}_t\|_1 \leq \frac{13m}{\rho_{\text{Stoc}}}$

641 and in the same way  $\sup_{t \in [T], x \in \mathcal{X}} \|u_t^{\mathbb{P}}(x)\|_1 \leq 2 \frac{13m}{\rho_{\text{Stoc}}}$ .

642 Define  $\xi$  as the best strategy that satisfies the constraints, i.e.,  $\text{Opt}_{\text{Stoc}} := T \mathbb{E}_{x \sim \xi} [\bar{f}(x)]$  and  
 643  $\mathbb{E}_{x \sim \xi} [\bar{g}_i(x)] \leq 0$ . The no-regret property of  $\text{Alg}_{\mathbb{P}}$  with respect to  $\xi$  gives that with probability  $1 - \delta$   
 644 it holds:

$$\begin{aligned} & \sum_{t \in [T]} [f_t(x_t) - \langle \boldsymbol{\lambda}_t, \mathbf{g}_t(x_t) \rangle] \\ & \geq \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} [f_t(x) - \langle \boldsymbol{\lambda}_t, \mathbf{g}_t(x) \rangle] \right] - \left( 2 \frac{13m}{\rho_{\text{Stoc}}} \right)^2 \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) \\ & \geq \mathbb{E}_{x \sim \xi} \left[ \sum_{t \in [T]} [\bar{f}(x) - \langle \boldsymbol{\lambda}_t, \bar{\mathbf{g}}(x) \rangle] \right] - 676 \left( \frac{m}{\rho_{\text{Stoc}}} \right)^2 \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) - 2 \left( \frac{13m}{\rho_{\text{Stoc}}} \right) E_{T,\delta} \\ & = T \text{Opt}_{\text{Stoc}} - 676 \left( \frac{m}{\rho_{\text{Stoc}}} \right)^2 \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) - \frac{26m}{\rho_{\text{Stoc}}} E_{T,\delta}, \end{aligned}$$

645 where the second inequality follows from Lemma B.3 with  $M := \frac{13m}{\rho_{\text{Stoc}}}$ .

646 Moreover, the no-regret property of the dual regret minimizer  $\text{Alg}_{\mathbb{D}}$ , with respect to  $\boldsymbol{\lambda}^* = \mathbf{0}$ , gives  
 647 that:

$$\sum_{t \in [T]} \langle \boldsymbol{\lambda}^* - \boldsymbol{\lambda}_t, \mathbf{g}_t(x_t) \rangle \leq \frac{1}{2} \eta_{\text{OGD}} mT.$$

648 Finally, we can combine everything from which follows that:

$$\text{Rew} \geq \text{Opt}_{\text{Stoc}} - 676 \left( \frac{m}{\rho_{\text{Stoc}}} \right)^2 \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}) - \frac{26m}{\rho_{\text{Stoc}}} E_{T,\delta} - \frac{1}{2} \eta_{\text{OGD}} mT.$$

649 The proof is concluded by observing that  $\eta_{\text{OGD}} = \tilde{O}((m \overline{R}^{\mathbb{P}}_{T,\delta}(\mathcal{X}))^{-1})$  and  $E_{T,\delta} = \tilde{O}(\sqrt{T})$   $\square$

650 **D Proofs omitted from Section 8**

651 **Lemma 8.3.** *The error of  $\mathcal{O}_{\mathcal{L}}$  can be bounded as*

$$\text{Err}(\mathcal{O}_{\mathcal{L}}) \leq 2\text{Err}(\mathcal{O}_f) + 2 \left( \sup_{t \in [T]} \|\lambda_t\|_1 \right)^2 \sum_{i \in [m]} \text{Err}(\mathcal{O}_i).$$

652 *Proof.* Consider the following inequalities:

$$\begin{aligned} \text{Err}(\mathcal{O}_{\mathcal{L}}) &:= \sum_{t \in [T]} \left( \hat{\mathcal{L}}_t(z_t, a_t) - \bar{\mathcal{L}}(z_t, a_t) \right)^2 \\ &\leq 2 \sum_{t \in [T]} \left( \hat{f}_t(z_t, a_t) - \bar{f}(z_t, a_t) \right)^2 + 2 \sum_{t \in [T]} \left( \langle \lambda_t, \hat{\mathbf{g}}_t(z_t, a_t) \rangle - \langle \lambda_t, \bar{\mathbf{g}}(z_t, a_t) \rangle \right)^2 \\ &\quad \text{(By AM-GM inequality: } 2ab \leq a^2 + b^2 \text{ for } a, b \geq 0.) \\ &= 2 \cdot \text{Err}(\mathcal{O}_f) + 2 \sum_{t \in [T]} \left( \langle \lambda_t, \hat{\mathbf{g}}_t(z_t, a_t) - \bar{\mathbf{g}}(z_t, a_t) \rangle \right)^2 \\ &\leq 2 \cdot \text{Err}(\mathcal{O}_f) + 2 \sum_{t \in [T]} \|\lambda_t\|_1^2 \cdot \|\hat{\mathbf{g}}_t(z_t, a_t) - \bar{\mathbf{g}}(z_t, a_t)\|_{\infty}^2 \quad (\langle a, b \rangle \leq \|a\|_1 \cdot \|b\|_{\infty}) \\ &\leq 2 \cdot \text{Err}(\mathcal{O}_f) + 2 \left( \sup_{t \in [T]} \|\lambda_t\|_1 \right)^2 \cdot \sum_{t \in [T]} \|\hat{\mathbf{g}}_t(z_t, a_t) - \bar{\mathbf{g}}(z_t, a_t)\|_{\infty}^2 \\ &\leq 2 \cdot \text{Err}(\mathcal{O}_f) + 2 \left( \sup_{t \in [T]} \|\lambda_t\|_1 \right)^2 \cdot \sum_{t \in [T]} \sum_{i \in [m]} (\hat{g}_{t,i}(z_t, a_t) - \bar{g}_i(z_t, a_t))^2 \\ &= 2 \cdot \text{Err}(\mathcal{O}_f) + 2 \left( \sup_{t \in [T]} \|\lambda_t\|_1 \right)^2 \cdot \sum_{i \in [m]} \text{Err}(\mathcal{O}_i) \end{aligned}$$

653 which concludes the proof. □

654 **Lemma 8.4.** *Assume that  $\max\{\text{Err}(\mathcal{O}_f), \text{Err}(\mathcal{O}_i)\} \leq \overline{\text{Err}}$ . Then, we have that Algorithm 3 with*  
 655  $\eta_P := \sqrt{KT}$  *guarantees that  $\sup_{I=[t_1, t_2]} R_I^P(\Pi) = \tilde{O}\left(m \cdot \overline{\text{Err}} \cdot L^2 \cdot \sqrt{KT}\right)$  with high probability,*  
 656 *where  $L := \sup_{t \in [T], \pi \in \Pi} |u_t^P(\pi)|$ .*

657 *Proof.* Consider any interval  $I = [t_1, t_2] \subseteq [T]$ . Since the prediction error at each time  $t$  is positive,  
 658 one trivially has that:

$$\sum_{t \in [t_1, t_2]} \left( \hat{\mathcal{L}}_t(z_t, a_t) - \bar{\mathcal{L}}(z_t, a_t) \right)^2 \leq \text{Err}(\mathcal{O}_{\mathcal{L}}).$$

659 Then, applying Lemma 8.3 we have that:

$$\sum_{t \in [t_1, t_2]} \left( \hat{\mathcal{L}}_t(z_t, a_t) - \bar{\mathcal{L}}(z_t, a_t) \right)^2 \leq 2\text{Err}(\mathcal{O}_f) + 2 \sup_{t \in [T]} \|\lambda_t\|_1^2 \sum_{i \in [m]} \text{Err}(\mathcal{O}_i).$$

660 Moreover, by the assumption on the errors of the oracles it holds that:

$$\sum_{t \in [t_1, t_2]} \left( \hat{\mathcal{L}}_t(z_t, a_t) - \bar{\mathcal{L}}(z_t, a_t) \right)^2 \leq 2m(1 + \sup_{t \in [T]} \|\lambda_t\|_1^2) \overline{\text{Err}}. \quad (12)$$

661 Note that we could pretend that the algorithm starts at any time  $t_1 \in [T]$ , and the same analysis of  
 662 [25, Theorem 1] would hold, as their algorithm behavior does not depend on its past behavior. Hence,

663 the following holds:

$$\begin{aligned}
R_{\llbracket t_1, t_2 \rrbracket}^{\mathbb{P}}(\Pi) &:= \sup_{\pi \in \Pi} \sum_{t \in \llbracket t_1, t_2 \rrbracket} [u_t^{\mathbb{P}}(\pi) - u_t^{\mathbb{P}}(\pi_t)] \\
&:= \sup_{\pi \in \Pi} \sum_{t \in \llbracket t_1, t_2 \rrbracket} [\mathcal{L}_t(\pi(z_t)) - \mathcal{L}_t(\pi_t(z_t))] \\
&= \sup_{\pi \in \Pi} \sum_{t \in \llbracket t_1, t_2 \rrbracket} [\mathcal{L}_t(\pi(z_t)) - \mathcal{L}_t(a_t)] \\
&\leq \frac{\eta_{\mathbb{P}}}{2} \text{Err}(\mathcal{O}_{\mathcal{L}}) + 4\eta_{\mathbb{P}} \log\left(\frac{2T^2}{\delta}\right) + 2K \frac{T}{\eta_{\mathbb{P}}} + \sqrt{2T \log\left(\frac{2T^2}{\delta}\right)}
\end{aligned}$$

664 which holds with probability  $1 - \delta/(T^2)$ .

665 Thus, by an union bound, and combining it with Equation (12) we obtain that:

$$R_{\llbracket t_1, t_2 \rrbracket}^{\mathbb{P}}(\Pi) \leq \eta_{\mathbb{P}} m (1 + \sup_{t \in \llbracket T \rrbracket} \|\boldsymbol{\lambda}_t\|_1^2) \overline{\text{Err}} + 4\eta_{\mathbb{P}} \log\left(\frac{2T^2}{\delta}\right) + 2K \frac{T}{\eta_{\mathbb{P}}} + \sqrt{2T \log\left(\frac{2T^2}{\delta}\right)},$$

666 which holds with probability  $1 - \delta/T^2$ . Finally, by tuning  $\eta_{\mathbb{P}} = \sqrt{KT}$  and applying an union bound  
667 on all the  $T^2$  possible intervals  $\llbracket t_1, t_2 \rrbracket$ , we obtain that with probability  $1 - \delta$  it holds that:

$$\sup_{I = \llbracket t_1, t_2 \rrbracket} R_{\llbracket t_1, t_2 \rrbracket}^{\mathbb{P}}(\Pi) \leq 504 \cdot m \overline{\text{Err}} L^2 \log(T^2/\delta) \sqrt{KT}.$$

668

□