

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

# Fast Algorithms for Online Stochastic Convex Programming

Shipra Agrawal

Department of Industrial Engineering and Operations Research, Columbia University, New York, NY 10027,  
shipra@ieor.columbia.edu

Nikhil R. Devanur

Microsoft Research, Redmond, WA. nikdev@microsoft.com

We introduce the *Online Stochastic Convex Programming (OSCP)* problem as a generalization of online linear programming, and a very natural model for sequential decision making under uncertainty. In this problem, stochastic inputs are revealed over time, a decision has to be made each time before observing the future inputs, and goal is to maximize a concave objective under convex feasibility constraints on the aggregate reward/cost vectors. This models many revenue management and resource allocation problems, particularly those under nonlinear resource consumption costs and risk-sensitive utility functions. We provide fast near-optimal algorithms for this problem under both i.i.d. and random order of arrival assumption on the inputs. Our algorithms are based on primal-dual paradigm, and we use online learning as a black-box to incrementally learn the optimal value of dual variables. Additionally, we demonstrate that under certain smoothness assumptions, online stochastic convex programming problem can be solved with much lower regret than its linear counterpart, thus providing incentive for considering smooth convex relaxations in practical settings. Even for online packing (a well studied special case of online linear programming), our techniques yield significantly faster and conceptually simpler algorithms than the state-of-the-art, with optimal guarantees.

*Key words:* convex programming; sequential decisions; online learning; Fenchel dual; display advertising; revenue management

*History:* A preliminary version of this paper appeared in Proceedings of SIAM Symposium on Discrete Algorithms (SODA) 2015.

## 1. Introduction

The theory of online matching and its generalizations has been a great success story that has had a significant impact on practice. The problems considered in this area are largely motivated by online advertising, and the theory has influenced how real advertising systems are run. As an example, the algorithms given by [Devanur et al. \(2011a\)](#) were used at Microsoft, by the “delivery engine” that decides which display ads are shown on its “properties” such as webpages, Skype, Xbox, etc.

In one of the most basic problem formulations in online advertising, an “impression” can be allocated to one of many given advertisers, assigning an impression  $i$  to advertiser  $a$  generates a value  $v_{ai}$ , and an advertiser  $a$  can be allocated at most  $G_a$  impressions. The goal is to maximize the value of the allocation. In another variant, advertisers pay per click and have budget constraints on their total payment, instead of the capacity constraints as above. More sophisticated formulations consider the option to show multiple ads on one webpage, which means you can pick among various configurations of ads. Each configuration still provides some value which is to be maximized, and advertisers have either capacity or budget constraints.

While the online budgeted matching algorithms, like the one in [Devanur et al. \(2011a\)](#) (DJSW algorithm), are used in practice, the actual problem has some aspects that are not captured by the formulations considered there. For instance, the actual objective function is not just a linear function, such as the sum of the values. There is a penalty for “under-delivering” impressions to an advertiser that increases with the amount of under-delivery. This translates into an objective that is a concave function of the total number of impressions assigned to an advertiser. Another consideration is the diversity of the impressions assigned. An advertiser targeting a certain segment of the population expects a representative sample of the entire population ([Ghosh et al. 2009](#)). In order to avoid deviating from this ideal too much, there are certain (convex) penalty functions in the objective that punish such deviations. The ‘essentially linear’ formulations of online matching or online packing/covering considered in the literature cannot handle these extensions. In this paper, we introduce a very general online *convex programming* framework, which allows any (Lipschitz) concave objective function and convex constraints on average reward/cost vectors, and present provably optimal algorithms for it.

Following related work on online packing and covering problems, we consider two closely related stochastic input models, the random permutation and the i.i.d. model. In the random permutation model, an adversary picks the *set* of inputs, which are then presented to the algorithm in a random order. In the i.i.d. model, the adversary picks a distribution over inputs that is unknown to the algorithm, and the algorithm receives i.i.d. samples from this distribution. The random permutation model is stronger than the i.i.d. model, any algorithm that works for the random permutation model also works for the i.i.d. model. The difference between these two models is like the difference between sampling with and without replacement. This intuition says that the two models should be very similar to each other, but the DJSW algorithm was only known to work for the i.i.d. model, not for the random permutation model. Earlier algorithms by [Devanur and Hayes \(2009\)](#), [Agrawal et al. \(2014\)](#), [Feldman et al. \(2010a\)](#) worked for the random permutation model but gave worse guarantees. [Kesselheim et al. \(2014\)](#) gave an algorithm that matched the optimal guarantee of [Devanur et al. \(2011a\)](#) for the random permutation model, but suffered on efficiency of implementation, as we discuss in the following.

An important practical consideration in the design of online algorithms is that the time taken by the algorithm in a single step should be very small. For instance, the decision to allocate an impression must be made in “real-time”, in a matter of milliseconds. The DJSW algorithm of [Devanur et al. \(2011a\)](#) satisfies this requirement, but requires solving an LP ever so often, to estimate the value of an optimum solution. On the other hand, the algorithm in [Kesselheim et al. \(2014\)](#) requires solving an LP in every step, making it not practical. In this paper, we give an algorithm that only requires solving a single LP (for *online packing and covering* problems), making it even faster than the DJSW algorithm. This improvement comes from the fact that in our algorithm the error in the estimate of the optimal solution only occurs in the second order error bounds and hence we can tolerate much bigger errors in such an estimate.

To summarize the comparison to work on online packing and covering, the DJSW algorithm is fast and works for the i.i.d. model but not for the random permutation model. The algorithm by [Kesselheim et al. \(2014\)](#) works for the random permutation model but is slow. We get the best of both worlds, our algorithm is fast, and works for the random permutation model. Moreover, our proof formalizes the

intuition mentioned earlier that the difference between i.i.d and the random permutation models is like the difference between sampling with and without replacement.

Furthermore, for the first time, we provide a general, provably optimal approach for handling global *convex* constraints and objective in online decision making. Our primal-dual algorithmic techniques employ fundamental concepts of Fenchel duality in convex programming, and use online learning methods as a blackbox to learn the dual variables. Starting from [Mehta et al. \(2007\)](#), it was conjectured that there is some relation between these problems and online learning or the “experts” problem, but no formal connection was known. We show such a formal connection, and demonstrate that getting better guarantees for these problems boils down to getting better “low-regret” guarantees for certain online learning problems. Even for the special case of online packing and covering, this gives much simpler proofs and improved algorithms than the earlier work.

To summarize, our contributions are as follows.

1. We present algorithms with near-optimal guarantees for a very general online convex programming problem, in a stochastic setting.
2. Our algorithms are primal-dual algorithms that are fast and simple, and work for the random permutation model. Our proof techniques formalize the intuition that the random permutation and the i.i.d models are not very different.
3. We establish a formal connection between these problems and online learning.

### 1.1. Other Related Work

The seminal paper of [Mehta et al. \(2007\)](#) introduced the so called “Adwords” problem, motivated by the allocation of ad slots on search engines, and started a slew of research into generalizations of the online bipartite matching problem ([Karp et al. 1990](#)). For the worst-case model, the optimal competitive ratio is  $1 - 1/e$ , which can be achieved for a fairly general setting ([Buchbinder et al. 2007](#), [Aggarwal et al. 2011](#), [Feldman et al. 2009a](#), [Devanur et al. 2013](#)). A special case of an objective with a concave function was considered in [Devanur and Jain \(2012\)](#).

In order to circumvent the impossibility results in the traditional worst-case models, stochastic models such as the random permutation model and the i.i.d model were introduced ([Goel and Mehta 2008](#),

Devanur and Hayes 2009, Vee et al. 2010, Devanur et al. 2011a). The dominant theme for these stochastic models has been asymptotic guarantees, that show that the competitive ratio tends to 1 as the “bid-to-budget” ratio tends to 0 (as was first shown by Devanur and Hayes (2009)). The focus then is the *convergence rate*, the rate at which the competitive ratio tends to 1 as a function of the bid-to-budget ratio. Feldman et al. (2010a), Agrawal et al. (2014) gave improved convergence rates for the random permutation model and generalized the result to an *online packing* problem. Recently, Chen and Wang (2013) extended these ideas to the concave returns problem of Devanur and Jain (2012). Devanur et al. (2011a) gave the optimal convergence rate for the online packing problem in the closely related i.i.d. model. Kesselheim et al. (2014) matched these bounds for the random permutation model, and further improved the bounds either when the bid-to-budget ratio is large, or when the instances are sparse. This line of research has also had significant impact on the practice of ad allocation with most of the big ad allocation platforms using algorithms influenced by these papers (Feldman et al. 2010b, Karande et al. 2013, Chen et al. 2011, 2012, Chakrabarti and Vee 2012).

Some versions of these problems also appear in literature under the name of ‘secretary problems’. However the dominant theme in research on secretary problems is to aim for a constant competitive ratio while not making any assumption about “bid-to-budget” ratio (a notable exception is (Kleinberg 2005)).

Another interesting line of research has been for the case of bipartite matching. Feldman et al. (2009b), Bahmani and Kapralov (2010), Manshadi et al. (2011) gave algorithms with competitive ratios better than  $1 - 1/e$  for the known distribution case, and Karande et al. (2011), Mahdian and Yan (2011) did the same for the random permutation model. Other variations such as models for combining algorithms from worst-case and average case, and achieving simultaneous guarantees have also been studied (Mahdian et al. 2012, Mirrokni et al. 2012).

A closely related problem is called the “Bandits with Knapsacks” problem (Badanidiyuru et al. 2013), which is similar to the online stochastic packing problem. The bandit aspect is different: the algorithm picks an “arm” of the bandit at each time, and makes observations (cost, reward, etc.), which are i.i.d samples that depend on the arm. There is persistence in the available set of choices across time as the

arms are persistent. In the online packing problem, the set of options in one time step are unrelated to the other time steps. Due to this, the main aspect of the bandit problem, the explore-exploit trade off in estimating the expectations of the observations for all arms, is absent from the online packing problem.

In an earlier paper (Agrawal and Devanur 2014), we generalized Bandits with Knapsacks to include general convex constraints and concave rewards, which is analogous to our generalization of the online packing to online convex programming here. Our high level ideas of using Fenchel duality for ‘linearization’ and online learning algorithms for estimating the dual variables is inspired by the use of similar ideas in (Agrawal and Devanur 2014). Consequently, we obtain algorithms that are very similar looking to those in (Agrawal and Devanur 2014). There are some significant differences in the proof techniques, however, due to the differences in the two problems mentioned in the previous paragraph. Also, the analysis for the random permutation model, and our adaptations (for the online packing problem) to get competitive ratios instead of regret bounds, were entirely absent from (Agrawal and Devanur 2014).

The online packing problem is also closely related to the Blackwell approachability problem (Blackwell 1956). The use of online learning algorithms to solve the Blackwell approachability problem (Abernethy et al. 2011) is similar to our use of online learning algorithms.

Concurrently and independently, Gupta and Molinaro (2014) found results for online linear programming that are similar to some of ours: they also show how to get competitive ratio bounds for the online packing problem in the random permutation model via a connection to the experts problem. For the guarantees that hold “in expectation”, their bounds are the same as ours. For the guarantees that hold “with high probability”, they show bounds without an extra  $\sqrt{\log T}$  factor that we get. They do not consider the more general convex programming framework.

## 1.2. Organization:

Section 2 contains the problem and the input model definitions, and statement of the main results. Section 3 provides some background material on online learning and Fenchel duality. Section 4 illustrates the basic ideas using a special case with only convex feasibility constraints. Section 5 gives the algorithm, results and proof techniques for the general online stochastic convex programming. Section 6 gives

tighter bounds for the special case of the online packing problem. Section 7 provides a new algorithm with stronger regret bounds for smooth functions.

## 2. Problem definition and main results

The following problem captures a very general setting of online optimization problems with global constraints and utility functions.

DEFINITION 1. [ONLINE STOCHASTIC CONVEX PROGRAMMING (OSCP)] We receive as initial input, the description of a concave function  $f$  over a bounded domain  $\subseteq \mathbb{R}^d$ , which we may assume is  $[0, 1]^d$  w.l.o.g, and a convex set  $S \subseteq [0, 1]^d$ . Subsequently we proceed in steps, at every time step  $t = 1, \dots, T$ , we receive a set  $A_t \subseteq [0, 1]^d$  of  $d$ -dimensional vectors. We have to pick one vector  $\mathbf{v}_t^\dagger \in A_t$  before proceeding to time step  $t + 1$ , using only information until time  $t$ . Let  $\mathbf{v}_{\text{avg}}^\dagger := \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^\dagger$ . The goal is to

$$\text{maximize } f(\mathbf{v}_{\text{avg}}^\dagger) \text{ subject to } \mathbf{v}_{\text{avg}}^\dagger \in S.$$

We assume that the instance is *always feasible*, i.e., there is a choice of  $\mathbf{v}_t \in A_t \forall t$  such that  $\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t \in S$ .

To appreciate the generality of above formulation, observe that objectives of form  $f(\sum_t g_t(\mathbf{v}_t^\dagger))$  can also be formulated in above by simply replacing every vector  $\mathbf{v}_t \in A_t$  by  $g_t(\mathbf{v}_t)$ , as long as  $g_t(\cdot)$  is bounded. Similarly, constraints of form  $h(\sum_t g_t(\mathbf{v}_t^\dagger)) \leq 0$  can be encoded in above formation, for any convex function  $h$ . Furthermore, even though vectors in  $A_t$  are required to be in  $[0, 1]^d$  in the above definition, one can handle vectors in  $[-1, 1]^d$  by replacing every  $\mathbf{v}_t \in A_t$  by  $\frac{\mathbf{v}_t + 1}{2} \in [0, 1]^d$ , and replacing the objective  $f(\mathbf{v}_{\text{avg}}^\dagger)$  and the constraint  $\mathbf{v}_{\text{avg}}^\dagger \in S$  by  $f(2\mathbf{v}_{\text{avg}}^\dagger - 1)$  and  $2\mathbf{v}_{\text{avg}}^\dagger - 1 \in S$ , respectively. The latter is equivalent to  $\mathbf{v}_{\text{avg}}^\dagger \in S'$  for a suitably modified convex set  $S'$ .

### 2.1. Stochastic Input Models:

In the random permutation (RP) model, there are  $T$  sets  $X_1, \dots, X_T$  fixed in advance but unknown to the algorithm, and these come in a uniformly random order (given by a random permutation  $\pi$ ) as the sequence  $A_1 = X_{\pi(1)}, \dots, A_T = X_{\pi(T)}$ . The number of time steps  $T$  is given to the algorithm in advance.

In the i.i.d, unknown distribution (IID) model, there is a distribution  $\mathcal{D}$  over subsets of  $[0, 1]^d$ , and for each  $t$ ,  $A_t$  is an independent sample from  $\mathcal{D}$ . The distribution  $\mathcal{D}$  is unknown to the algorithm.

It is known that the RP model is stronger than the IID model. The IID model can be thought of as a distribution over RP instances and therefore any guarantee for the RP model also carries over to the IID model. Henceforth, we will consider the RP model by default, unless otherwise mentioned.

## 2.2. Benchmarks.

We measure the performance of an algorithm with respect to a benchmark. The benchmark for the RP model is the *optimal offline* solution, i.e. the choice  $\mathbf{v}_t^* \in A_t$  that maximizes the function  $f$  of the average of these vectors while making sure that the average lies in  $S$ . We denote the value of this solution as the benchmark, OPT. This is a deterministic value since it does not depend on the randomness in the input, which is in the order of arrival. For the IID model, the offline optimal actually depends on the randomness in the input, and OPT denotes the expected value of the offline optimal solution.

## 2.3. Performance Measures.

While the standard measure in competitive analysis of online algorithms is a multiplicative error w.r.t the benchmark, we mostly adopt a concept of additive error that is common in online learning, called the *regret*. Since we make no assumptions about  $f$ , it could even be negative, so an additive error is more appropriate. For certain special cases where multiplicative errors or competitive ratios are more natural or desirable, we discuss how our algorithms and analysis can be adapted to get such guarantees. We define the following two (average) regret measures, one for the objective and another for the constraint.<sup>1</sup> Let  $d(\mathbf{v}, S)$  denote the distance of the vector  $\mathbf{v}$  from the set  $S$ , w.r.t. a given norm  $\|\cdot\|$ .

$$\text{avg-regret}_1(T) = \text{OPT} - f(\mathbf{v}_{\text{avg}}^\dagger), \text{ and}$$

$$\text{avg-regret}_2(T) = d(\mathbf{v}_{\text{avg}}^\dagger, S).$$

## 2.4. Main Results.

We now state the most general result we prove in this paper.

**THEOREM 1.** *There is an algorithm (Algorithm 2) that achieves the following regret guarantees for the Online Stochastic Convex Programming problem, in the RP model.*

$$\begin{aligned} \mathbb{E}[\text{avg-regret}_1(T)] &= (Z + L) \cdot O\left(\sqrt{\frac{C}{T}}\right) \\ \mathbb{E}[\text{avg-regret}_2(T)] &= O\left(\sqrt{\frac{C}{T}}\right) \end{aligned}$$



Here, the Big-Oh notation is hiding only universal constants.  $C$  depends on the norm  $\|\cdot\|$  used for defining distance. For Euclidean norm,  $C = d \log(d)$ . For  $L_\infty$  norm,  $C = \log(d)$ . The parameter  $Z$  captures the tradeoff between objective and constraints for the problem, its value is problem-dependent and is discussed in detail later in the text.  $L$  is the Lipschitz constant for  $f$  w.r.t. the same norm  $\|\cdot\|$  as used to measure the distance.

In the main text we provide more detailed result statements, which will also make clear the dependence of our regret bounds on the regret bounds available for online learning, and implications of using different norms. These regret bounds can also be converted to *high probability* results, with an additional  $\sqrt{\log T}$  factor in the regret. This extra factor comes from simply taking a union bound over all time steps. A more careful analysis could possibly get rid of this extra factor, as was shown in Gupta and Molinaro (2014) in case of online linear programming. These bounds are optimal, and this follows easily from an easy modification of a lower bound given by Agrawal et al. (2014) for the online packing problem.

We also consider the following interesting special cases.

*Feasibility problem:* In this case, there is no objective function  $f$ , and there is only the constraint given by the set  $S$ . The goal is to make sure that the average of the chosen vectors lies as close to  $S$  as possible, i.e., minimize  $d(\mathbf{v}_{\text{avg}}^\dagger, S)$ .

*Linear objective:* In this case, we assume that each vector  $\mathbf{v} \in A_t$  has an associated reward  $r \in [0, 1]$ . The objective is to maximize the total reward while making sure that the average of the vectors lies in  $S$ . This can be thought of as the special case where the vector you get is  $(\mathbf{v}, r)$ , and the constraint is only on the subspace defined by all coordinates of this vector except the last, while the objective is just the sum (or linear function) of its last coordinates.

*Online Packing/Covering LPs:* This is a well studied special case of linear objective. The packing constraints  $\sum_t \mathbf{v}_t^\dagger \leq B\mathbf{1}$  are equivalent to using constraint set  $S$  of the form  $\{\mathbf{v} : 0 \leq \mathbf{v} \leq \frac{B}{T}\mathbf{1}\}$ , where  $\mathbf{1}$  is the vector of all 1s and  $B > 0$  is some scalar. In this case, we also assume that the sets  $A_t$  always contain the origin, which corresponds to the option of “doing nothing”. The covering constraints are obtained when  $S$  is  $\{\mathbf{v} : \mathbf{v} \geq \frac{B}{T}\mathbf{1}\}$ .

For online packing, we provide the following tighter guarantee in terms of competitive ratio.

**THEOREM 2.** *For online stochastic packing problem, Algorithm 4 achieves a competitive ratio of  $1 - O(\epsilon)$  in the RP model, given any  $\epsilon > 0$  such that  $\min\{B, TOPT\} \geq \log(d)/\epsilon^2$ . Further, the algorithm has fast per-step updates, and needs to solve a sample LP at most once.*

### 3. Preliminaries

In this section, we provide background on some fundamental technical concepts and facts used in this paper.

#### 3.1. Fenchel duality.

As mentioned earlier, our algorithms are primal-dual algorithms. For the online packing problem, the LP duality framework (which is very well understood) is sufficient but for general convex programs we need the stronger framework of Fenchel duality. Below we provide some background on this useful mathematical concept.

Let  $h$  be a convex function defined on  $[0, 1]^d$ . We define  $h^*$  as the Fenchel conjugate of  $h$ ,

$$h^*(\boldsymbol{\theta}) := \max\{\mathbf{y} \cdot \boldsymbol{\theta} - h(\mathbf{y}) : \mathbf{y} \in [0, 1]^d\}.$$

Similarly for a concave function  $f$  on  $[0, 1]^d$ , define  $f^*(\boldsymbol{\theta}) := \max_{\mathbf{y} \in [0, 1]^d} \{\mathbf{y} \cdot \boldsymbol{\theta} + f(\mathbf{y})\}$ . Note that the Fenchel conjugates  $h^*$  and  $f^*$  are both convex functions of  $\boldsymbol{\theta}$ .

Suppose that at every point  $\mathbf{y}$ , every supergradient  $\mathbf{g}_{\mathbf{y}}$  of  $h$  (and  $f$ ) have bounded dual norm  $\|\mathbf{g}_{\mathbf{y}}\|_* \leq L$ . Then, the following dual relationship is known between  $h$  and  $h^*$  ( $f$  and  $f^*$ ).

**LEMMA 1.**  $h(\mathbf{z}) = \max_{\|\boldsymbol{\theta}\|_* \leq L} \{\boldsymbol{\theta} \cdot \mathbf{z} - h^*(\boldsymbol{\theta})\}$ ,  $f(\mathbf{z}) = \min_{\|\boldsymbol{\theta}\|_* \leq L} \{f^*(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{z}\}$

A special case is when  $h(\mathbf{y}) = d(\mathbf{y}, S)$  for some convex set  $S$ . This function is 1-Lipschitz with respect to norm  $\|\cdot\|$  used in the definition of distance. In this case,  $h^*(\boldsymbol{\theta}) = h_S(\boldsymbol{\theta}) := \max_{\mathbf{y} \in S} \boldsymbol{\theta} \cdot \mathbf{y}$ , and Lemma 1 specializes to the following relation (which also appears in [Abernethy et al. \(2011\)](#)).

$$d(\mathbf{y}, S) = \max\{\boldsymbol{\theta} \cdot \mathbf{y} - h_S(\boldsymbol{\theta}) : \|\boldsymbol{\theta}\|_* \leq 1\}. \quad (1)$$

### 3.2. Strong convexity/Smoothness Duality.

We first define strong convexity and smoothness.

DEFINITION 2. A function  $h : \mathcal{X} \rightarrow \mathbb{R}$  is  $\beta$ -strongly convex w.r.t. a norm  $\|\cdot\|$  if  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}, \mathbf{z} \in \partial h(\mathbf{x})$ ,

$$h(\mathbf{y}) - h(\mathbf{x}) \geq \mathbf{z} \cdot (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Equivalently for any  $\mathbf{x}, \mathbf{y}$  in the interior of  $\mathcal{X}$ , and all  $\alpha \in (0, 1)$ , we have that

$$\begin{aligned} h(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) &\geq \alpha h(\mathbf{x}) + (1 - \alpha) h(\mathbf{y}) \\ &\quad - \frac{\beta}{2} \alpha (1 - \alpha) \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

A function  $h$  is  $\beta$ -strongly concave if and only  $(-h)$  is  $\beta$ -strongly convex.

DEFINITION 3. A function  $h : \mathcal{X} \rightarrow \mathbb{R}$  is  $\beta$ -strongly smooth w.r.t. a norm  $\|\cdot\|$  if  $h$  is everywhere differentiable, and for all  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ , we have

$$\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}, |h(\mathbf{y}) - h(\mathbf{x}) - \nabla h(\mathbf{x}) \cdot (\mathbf{y} - \mathbf{x})| \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

The following lemma can be derived from the proof of Theorem 6 in [Kakade et al. \(2009\)](#). A proof is given in Appendix [B.1](#) for completeness.

LEMMA 2. *If  $h$  is convex and  $\beta$ -strongly smooth with respect to norm  $\|\cdot\|$ , then  $h^*(\boldsymbol{\theta}) = \max_{\mathbf{x} \in [0, 1]^d} \{\boldsymbol{\theta} \cdot \mathbf{x} - h(\mathbf{x})\}$  is  $\frac{1}{\beta}$ -strongly convex with respect to norm  $\|\cdot\|_*$  on domain  $\nabla_h = \{\nabla h(\mathbf{x}) : \mathbf{x} \in [0, 1]^d\}$ .*

### 3.3. Online Learning.

A well studied problem in online learning, called the Online Convex Optimization (OCO) problem, considers a  $T$  round game played between a learner and an adversary (nature), where at round  $t$ , the player chooses a  $\boldsymbol{\theta}_t \in W$ , and then the adversary picks a concave function  $g_t(\boldsymbol{\theta}_t) : W \rightarrow \mathbb{R}$ . The player's choice  $\boldsymbol{\theta}_t$  may only depend on the adversary's choices in the previous rounds. The goal of the player is to minimize regret defined as the difference between the player's objective value and the value of the best single choice in hindsight:

$$\mathcal{R}(T) := \max_{\boldsymbol{\theta} \in W} \sum_{t=1}^T g_t(\boldsymbol{\theta}) - \sum_{t=1}^T g_t(\boldsymbol{\theta}_t)$$

Some popular algorithms for OCO are online mirror descent (OMD) algorithm and online gradient descent, which have very fast per step update rules, and provide the following regret guarantees. More details about these algorithms and their regret guarantees are in Appendix B.2.

LEMMA 3. *Shalev-Shwartz (2012)* There is an algorithm for the OCO problem that achieves regret

$$\mathcal{R}(T) = O(G\sqrt{DT}),$$

where  $D$  is the diameter of  $W$  and  $G$  is an upper bound on the norm of gradient of  $g_t(\boldsymbol{\theta})$  for all  $t$ . The value of these parameters are problem specific.

In particular, following corollary can be derived, which will be useful for our purpose. Details are in Appendix B.2.

COROLLARY 1. For  $g_t(\boldsymbol{\theta})$  of form  $g_t(\boldsymbol{\theta}) = \boldsymbol{\theta} \cdot \mathbf{z} - h^*(\boldsymbol{\theta})$  and  $W = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_* \leq L\}$ , where  $h$  is an  $L$ -Lipschitz function, OCO algorithms achieve regret bounds of  $\mathcal{R}(T) \leq O(L\sqrt{dT})$  for Euclidian norm, and  $O(L\sqrt{\log(d)T})$  for  $L_\infty$ .

For optimization over a simplex ( $\|\boldsymbol{\theta}\|_1 = 1, \boldsymbol{\theta} \geq 0$ ), the *multiplicative weight update* algorithm (generalization by Arora et al. (2012)) is very fast and efficient. It is conventional to state this algorithm in terms of learning distribution over  $d$  experts, with  $\theta_j$  denoting the probability of  $j^{\text{th}}$  expert, and  $g_{t,j} = g_t(\mathbf{e}_j)$  being the outcome of  $j^{\text{th}}$  expert. To handle domain  $\|\boldsymbol{\theta}\|_1 \leq 1, \boldsymbol{\theta} \geq 0$ , and outcomes given by  $g_t(\boldsymbol{\theta})$ , we consider the problem with  $d + 1$  experts, with  $g_{t,j} = g_t(\mathbf{e}_j)$  being the outcome for expert  $j$  at time step  $t$ , and the outcome of  $d + 1$  expert always being  $g_t(\mathbf{0}) = 0$ . Then, the step  $t$  update of this algorithm takes the following form, given that  $-R \leq g_{t,j} \leq M$  and a parameter  $\epsilon > 0$ ,

$$\boldsymbol{\theta}_{t+1,j} = \frac{w_{t,j}}{1 + \sum_j w_{t,j}}, \text{ where } w_{t,j} = \begin{cases} w_{t-1,j}(1 + \epsilon)^{g_{t,j}/M} & \text{if } g_{t,j} > 0, \\ w_{t-1,j}(1 - \epsilon)^{-g_{t,j}/M} & \text{if } g_{t,j} \leq 0 \end{cases} \quad (2)$$

The multiplicative weight update algorithm then provides the following guarantees for  $\boldsymbol{\theta}_t$ s generated in the above manner.

LEMMA 4. *Arora et al. (2012)* Let  $-R \leq g_{t,j} \leq M, j = 1, \dots, d$  denote the outcomes of  $d$  experts at time  $t$  and expert  $d + 1$  always generates outcome 0. Then, for all  $0 < \epsilon \leq \frac{1}{2}$ , for  $\boldsymbol{\theta}_t$ s generated by the multiplicative weight update algorithm,

$$\sum_{t=1}^T \boldsymbol{\theta}_t \cdot \mathbf{g}_t \geq (1 - \epsilon) \left( \sum_{\geq 0} g_{t,j} \right) + (1 + \epsilon) \left( \sum_{< 0} g_{t,j} \right) - \frac{M \ln(d+1)}{\epsilon}, \text{ and,}$$

$$\sum_{t=1}^T \boldsymbol{\theta}_t \cdot \mathbf{g}_t \geq -\frac{M \ln(d+1)}{\epsilon},$$

where  $\geq 0$  and  $< 0$  refer to the rounds  $t$  where  $g_{t,j} \geq 0$  and  $g_{t,j} < 0$  respectively.

By concavity of  $g_t(\boldsymbol{\theta})$ , we obtain the following corollary.

**COROLLARY 2.** For domain  $W = \{\|\boldsymbol{\theta}\|_1 \leq 1, \boldsymbol{\theta} \geq 0\}$ , given that  $-R \leq g_t(\boldsymbol{\theta}_t) \leq M$ , and for all  $0 < \epsilon \leq \frac{1}{2}$ , using the multiplicative weight update algorithm we obtain that for any  $j$ ,

$$\sum_{t=1}^T g_t(\boldsymbol{\theta}_t) \geq \max \left\{ (1 - \epsilon) \left( \sum_{\geq 0} g_t(\mathbf{e}_j) \right) + (1 + \epsilon) \left( \sum_{< 0} g_t(\mathbf{e}_j) \right) - \frac{M \ln(d+1)}{\epsilon}, -\frac{M \ln(d+1)}{\epsilon} \right\}.$$

For strongly concave functions, even stronger *logarithmic* regret bounds can be achieved.

**LEMMA 5.** *Hazan et al. (2007)* Suppose that  $g_t$  is  $H$ -strongly concave for all  $t$ , and  $G \geq 0$  is an upper bound on the norm of the gradient, i.e.  $\|\nabla g_t(\boldsymbol{\theta})\| \leq G$ , for all  $t$ . Then the online gradient descent algorithm achieves the following guarantees for OCO: for all  $T \geq 1$ ,

$$\mathcal{R}(T) \leq \frac{G^2}{H} \log(T).$$

## 4. Feasibility Problem

It will be useful to first illustrate our algorithm and proof techniques for the special case of the feasibility problem. In this special case of online stochastic CP, there is no objective function  $f$ , and the aim of the algorithm is to have  $\mathbf{v}_{\text{avg}}^\dagger$  be in the set  $S$ . The performance of the algorithm is measured by the distance from the set  $S$ , i.e.,  $d(\mathbf{v}_{\text{avg}}^\dagger, S)$ . We assume that the instance is *always feasible*, i.e., there exist  $\mathbf{v}_t^* \in A_t \forall t$  such that  $\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^* \in S$ .

The basic idea behind our algorithm is as follows. Suppose that instead of minimizing a convex function such as  $d(\mathbf{v}_{\text{avg}}^\dagger, S)$  we had to minimize a linear function such as  $\boldsymbol{\theta} \cdot \mathbf{v}_{\text{avg}}^\dagger$ . This would be extremely easy since the problem then separates into small subproblems where at each time step we can simply solve  $\min_{\mathbf{v}_t^\dagger \in A_t} \boldsymbol{\theta} \cdot \mathbf{v}_t^\dagger$ . In fact, convex programming duality guarantees exactly this – that there is a  $\boldsymbol{\theta}^*$ , such that an optimal (i.e., feasible) solution is  $\mathbf{v}_t^* = \arg \min_{\mathbf{v} \in A_t} \boldsymbol{\theta}^* \cdot \mathbf{v}$ , however, we don't know  $\boldsymbol{\theta}^*$ . This is where online learning comes into play. Online learning algorithms can provide a  $\boldsymbol{\theta}_t$  at every time  $t$  using only the observations before time  $t$ , which together provide a good approximation to the best  $\boldsymbol{\theta}$  in hindsight.

**Algorithm 1** Feasibility problemInitialize  $\boldsymbol{\theta}_1$ .**for all**  $t = 1, \dots, T$  **do**Set  $\mathbf{v}_t^\dagger = \arg \min_{\mathbf{v} \in A_t} \boldsymbol{\theta}_t \cdot \mathbf{v}$ Choose  $\boldsymbol{\theta}_{t+1}$  by doing an OCO update with  $g_t(\boldsymbol{\theta}) = \boldsymbol{\theta} \cdot \mathbf{v}_t^\dagger - h_S(\boldsymbol{\theta})$ , and domain  $W = \{\|\boldsymbol{\theta}\|_* \leq 1\}$ .**end for**

Here  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ , the norm used in the distance function. The updates required for selecting  $\boldsymbol{\theta}_{t+1}$ , given  $\boldsymbol{\theta}_t$  and  $g_t(\cdot)$ , are given as Equation 13 and Equation 2 for OMD and multiplicative weight update algorithm, respectively. As discussed there, these updates are simple and fast, and do not require solving any complex optimization problems.

**THEOREM 3.** *Algorithm 1 achieves the following regret bound for the Feasibility Problem in the RP model of stochastic inputs:*

$$\mathbb{E}[\text{avg-regret}_2(T)] := \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \leq O\left(\frac{\mathcal{R}(T)}{T} + \|\mathbf{1}_d\| \sqrt{\frac{s \log(d)}{T}}\right).$$

where  $\mathcal{R}(T)$  denotes the regret for OCO with functions  $g_t(\boldsymbol{\theta})$  and domain  $W$ , as defined in Section 3.3. And,  $s \leq 1$  is the coordinate-wise largest value a vector in  $S$  can take. The parameter  $s$  can be used to obtain tighter problem-specific bounds.

*Proof.* From Fenchel duality, and by OCO guarantees,

$$\begin{aligned} d(\mathbf{v}_{\text{avg}}^\dagger, S) &= \max_{\|\boldsymbol{\theta}\|_* \leq 1} \boldsymbol{\theta} \cdot \mathbf{v}_{\text{avg}}^\dagger - h_S(\boldsymbol{\theta}) \\ &= \max_{\|\boldsymbol{\theta}\|_* \leq 1} \frac{1}{T} \sum_t g_t(\boldsymbol{\theta}) \\ &\leq \frac{1}{T} \sum_t g_t(\boldsymbol{\theta}_t) + \frac{1}{T} \mathcal{R}(T). \end{aligned}$$

In Lemma 6, we upper bound  $\mathbb{E}[\frac{1}{T} \sum_t g_t(\boldsymbol{\theta}_t)]$  to obtain the statement of the theorem.  $\square$

**LEMMA 6.**  $\mathbb{E}[\sum_t g_t(\boldsymbol{\theta}_t)] \leq O(\|\mathbf{1}_d\| \sqrt{s \log(d) T})$ , where  $s = \max_{\mathbf{v} \in S} \max_j v_j \leq 1$ , and  $\|\cdot\|$  is the norm used in the distance function.

*Proof.* Let  $\mathcal{F}_{t-1}$  denote the observations and decisions until time  $t-1$ . Note that  $\boldsymbol{\theta}_t$  is completely determined by  $\mathcal{F}_{t-1}$ . Let  $\mathbf{v}_{X_t}$  denote the option chosen to satisfy request  $X_t$  by the offline optimal

(feasible) solution, and let  $\mathbf{v}_t^* = \mathbf{v}_{A_t}$ . Then, since  $A_t = X_s$ , for  $s = 1, \dots, T$  with equal probability, we have that  $\mathbb{E}[\mathbf{v}_t^*] = \frac{1}{T}(\mathbf{v}_{X_1} + \dots + \mathbf{v}_{X_T}) \in S$ . Therefore, due to the manner in which  $\mathbf{v}_t^\dagger$  was chosen by the algorithm, we have that

$$\begin{aligned} \mathbb{E}[g_t(\boldsymbol{\theta}_t)|\mathcal{F}_{t-1}] &= \mathbb{E}[\boldsymbol{\theta}_t \cdot \mathbf{v}_t^\dagger - h_S(\boldsymbol{\theta}_t)|\mathcal{F}_{t-1}] \\ &\leq \mathbb{E}[\boldsymbol{\theta}_t \cdot \mathbf{v}_t^* - h_S(\boldsymbol{\theta}_t)|\mathcal{F}_{t-1}] \\ &= \boldsymbol{\theta}_t \cdot \mathbb{E}[\mathbf{v}_t^*] - h_S(\boldsymbol{\theta}_t) + \boldsymbol{\theta}_t \cdot (\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]). \end{aligned}$$

Now, by the Fenchel dual representation of distance, for any  $\mathbf{v}, \boldsymbol{\theta}'$  such that  $\|\boldsymbol{\theta}'\|_* \leq 1$ ,  $d(\mathbf{v}, S) = \max_{\|\boldsymbol{\theta}'\|_* \leq 1} \boldsymbol{\theta}' \cdot \mathbf{v} - h_S(\boldsymbol{\theta}') \geq \boldsymbol{\theta}' \cdot \mathbf{v} - h_S(\boldsymbol{\theta}')$ . Using this observation along with  $\mathbb{E}[\mathbf{v}_t^*] \in S$ , we obtain from above,

$$\begin{aligned} \mathbb{E}[g_t(\boldsymbol{\theta}_t)|\mathcal{F}_{t-1}] &\leq d(\mathbb{E}[\mathbf{v}_t^*], S) + \boldsymbol{\theta}_t \cdot (\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]) \\ &= 0 + \boldsymbol{\theta}_t \cdot (\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]) \\ &\leq \|\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\|, \end{aligned} \tag{3}$$

where the last inequality used the condition  $\|\boldsymbol{\theta}_t\|_* \leq 1$ .

Note that under independence assumption (IID model), we would have  $\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] = \mathbb{E}[\mathbf{v}_t^*]$ , so that the above inequality would suffice to give the required bound. However, in random permutation (RP) model, the observations till time  $t-1$  restrict the set of possible permutations. Conditional on realization  $A_1 = X_{\pi(1)}, \dots, A_{t-1} = X_{\pi(t-1)}$  until time  $t-1$ , for a given ordering  $\pi$ , we have that  $A_t$  is one of the *remaining sets* with equal probability. So,  $\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] = \frac{1}{T-t+1}(\mathbf{v}_{X_{\pi(t)}} + \dots + \mathbf{v}_{X_{\pi(T)}})$ , for any ordering  $\pi$  that agrees with  $\mathcal{F}_{t-1}$  on the first  $t-1$  indices.

Next, we bound the gap  $\|\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\|$  under random permutation assumption. For any given ordering  $\pi$ , define  $\mathbf{w}_{t,\pi} = \frac{\mathbf{v}_{X_{\pi(1)}} + \dots + \mathbf{v}_{X_{\pi(t)}}}{t}$ . Also, for given ordering  $\pi$ , define  $\pi'$  as the reverse ordering. Then,  $\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] = \mathbf{w}_{T-t+1,\pi'}$ , for any ordering  $\pi$  that agrees with  $\mathcal{F}_{t-1}$  on the first  $t-1$  indices. Now, the input ordering  $\pi$  observed by the algorithm agrees with all the filtrations  $\mathcal{F}_1, \dots, \mathcal{F}_{T-1}$ , and therefore taking  $\pi'$  as the reverse of this ordering, we have that

$$\begin{aligned} \sum_{t=1}^T \|\mathbb{E}[\mathbf{v}_t^*|\mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\| &= \sum_{t=1}^T \|\mathbf{w}_{T-t+1,\pi'} - \mathbb{E}[\mathbf{v}_t^*]\| \\ &= \sum_{t=1}^T \|\mathbf{w}_{t,\pi'} - \mathbb{E}[\mathbf{v}_t^*]\|. \end{aligned}$$

Due to the random permutation assumption, the input ordering  $\pi$ , and hence the reverse ordering  $\pi'$  in above, is a uniformly random permutation. Also, taking expectation over uniformly random permutations  $\sigma$ ,  $\mathbb{E}[\mathbf{w}_{t,\sigma}] = \frac{\mathbf{v}_{X_1} + \dots + \mathbf{v}_{X_T}}{T} = \mathbb{E}[\mathbf{v}_t^*]$ . And, therefore,

$$\sum_{t=1}^T \|\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\| = \sum_{t=1}^T \|\mathbf{w}_{t,\pi} - \mathbb{E}[\mathbf{w}_{t,\sigma}]\|, \quad (4)$$

where  $\pi$  is a uniformly random permutation. Taking outer expectations, and using (3), this implies,

$$\begin{aligned} \mathbb{E}\left[\sum_t g_t(\boldsymbol{\theta}_t)\right] &\leq \mathbb{E}\left[\sum_t \|\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\|\right] \\ &= \mathbb{E}\left[\sum_t \|\mathbf{w}_{t,\pi} - \mathbb{E}[\mathbf{w}_{t,\sigma}]\|\right]. \end{aligned}$$

Observe that for uniformly random permutation  $\pi$ ,  $\mathbf{w}_{t,\pi}$  can be viewed as the average of  $t$  vectors sampled uniformly *without replacement* from the ground set  $\{\mathbf{v}_{X_1}, \dots, \mathbf{v}_{X_T}\}$  of  $T$  vectors. We use Chernoff-Hoeffding type concentration bounds for sampling without replacement (refer to Appendix C for details), to obtain,

$$\mathbb{E}[\|\mathbf{w}_{t,\pi} - \mathbb{E}[\mathbf{w}_{t,\sigma}]\|] \leq O(\|\mathbf{1}_d\| \sqrt{\frac{s \log(d)}{t}}). \quad (5)$$

The lemma statement then follows by summing up these bounds over all  $t$ .  $\square$

REMARK 1. [*RP vs. IID*] For the IID model, since  $\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] = \mathbb{E}[\mathbf{v}_t^*]$ , we would get  $\sum_t \mathbb{E}[g_t(\boldsymbol{\theta}_t)] \leq 0$  directly from Equation (3). Thus, the quantity  $\mathbb{E}[\sum_t \|\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\|] \leq O(\|\mathbf{1}_d\| \sqrt{sT \log(d)})$  characterizes the gap between IID and RP models.

REMARK 2. [*High probability bounds*] The above analysis can be extended to bound the sum of *conditional expectations*  $\sum_t \mathbb{E}[g_t(\boldsymbol{\theta}_t) | \mathcal{F}_{t-1}] \leq \sum_t \|\mathbf{w}_{t,\pi} - \mathbb{E}[\mathbf{w}_{t,\sigma}]\|$  by  $O(\|\mathbf{1}_d\| \sqrt{T \log(dT/\rho)})$  with high probability  $1 - \rho$ . As a result, we obtain a high probability regret bound of  $O(\|\mathbf{1}_d\| \sqrt{\frac{\log(Td)}{T}})$ . Details are in Appendix C. For the IID model, this sum of conditional expectations is bounded by 0, so the resulting high probability bounds are slightly stronger, with no extra  $\sqrt{\log(T)}$  factor.

## 5. Online stochastic convex programming

In this section, we extend the algorithm from previous section to the general Online Stochastic Convex Programming (OSCP) problem, as defined in Section 2. Recall that the aim here is to maximize  $f(\mathbf{v}_{\text{avg}}^\dagger)$  while ensuring  $\mathbf{v}_{\text{avg}}^\dagger \in S$ .



A direct way to extend the algorithm from the previous section would be to reduce the convex program to the feasibility problem with constraint set  $S' = \{\mathbf{v} : f(\mathbf{v}) \geq \text{OPT}, \mathbf{v} \in S\}$ . However, this requires the knowledge of  $\text{OPT}$ . If  $\text{OPT}$  is estimated, the errors in the estimation of  $\text{OPT}$  at all time steps  $t$  would add up to the regret, thus this approach would tolerate very small  $\tilde{O}(\frac{1}{\sqrt{t}})$  per step estimation errors. In this section, we propose an alternate approach of combining objective value and distance from constraints using a parameter  $Z$ , which will capture the tradeoff between the two quantities. We may still need to estimate this parameter  $Z$ , however,  $Z$  will appear only in the second order regret terms, so that a constant factor approximation of  $Z$  will suffice to obtain optimal order of regret bounds. This makes the estimation task relatively easy and enable us to get better problem specific bounds. As a specific example, for the online packing problem, we can use  $Z = \frac{\text{OPT}}{(B/T)}$  so this approach requires only a constant factor approximation of  $\text{OPT}$  and the resulting algorithm obtains the optimal competitive ratio. (See Section 6 for more details.)

To illustrate the main ideas in our algorithm, let us start with the following assumption.

ASSUMPTION 1. Let  $\text{OPT}^\delta$  denote the optimal value of the offline problem that maximizes  $f(\frac{1}{T} \sum_t \mathbf{v}_t)$  with feasibility constraint relaxed to  $d(\frac{1}{T} \sum_t \mathbf{v}_t, S) \leq \delta$ . We are given a  $Z \geq 0$  such that that for all  $\delta \geq 0$ ,

$$\text{OPT}^\delta \leq \text{OPT} + Z\delta. \tag{6}$$

In fact, such a  $Z$  always exists, as shown by the following lemma.

LEMMA 7.  $\text{OPT}^\delta$  is a non-decreasing concave function of the constraint violation  $\delta$ , and its gradient at  $\delta = 0$  is the minimum value of  $Z$  that satisfies the property (6). This gradient is also equal to the value of the optimal dual variable corresponding to the distance constraint.

The proof of this lemma is provided in Appendix D. This fact is known for linear programs.

Below, we present an algorithm (Algorithm 2) for OSCP assuming we are given parameter  $Z$  as in Assumption 1. This algorithm is based on the same basic ideas as the algorithm for the feasibility problem in the previous section. Here, we linearize both objective and constraints using Fenchel duality, and estimate the corresponding dual variables using online learning as blackbox. And, we use parameter  $Z$  to combine objective with constraints. The resulting algorithm has very efficient per-step updates

and does not require solving a (sample) convex program in any step, and we prove that it achieves the regret bound stated in Theorem 1.

The regret of this algorithm (as stated in Theorem 1) scales with the value of  $Z$ , and it is desirable to use as small a value of  $Z$  as possible. If such a  $Z$  is not known, in Appendix F we demonstrate how we can approximate the optimal value of  $Z$  up to a constant factor by solving a logarithmic number of sample convex programs overall.

---

**Algorithm 2** Online convex programming
 

---

Initialize  $\theta_1, \phi_1$ .

**for all**  $t = 1, \dots, T$  **do**

    Choose option

$$\mathbf{v}_t^\dagger = \arg \max_{\mathbf{v} \in A_t} -\phi_t \cdot \mathbf{v} - 2(Z + L)\theta_t \cdot \mathbf{v}.$$

    Choose  $\theta_{t+1}$  by doing an OCO update for  $g_t(\theta) = \theta \cdot \mathbf{v}_t^\dagger - h_S(\theta)$  over domain  $W = \{\|\theta\|_* \leq 1\}$ .

    Choose  $\phi_{t+1}$  by doing an OCO update for  $\psi_t(\phi) = \phi \cdot \mathbf{v}_t^\dagger - (-f)^*(\phi)$  over domain  $U = \{\|\phi\|_* \leq L\}$ .

**end for**

---

A complete proof of Theorem 1, along with a more detailed theorem statement, is provided in Appendix E. Here, we provide the proof for the simpler case of *linear objective* discussed in Section 2. In this setting, each option in  $A_t$  is associated with a reward  $r$  in addition to the vector  $\mathbf{v}$ . And, at every time step  $t$ , the player chooses  $(r_t^\dagger, \mathbf{v}_t^\dagger)$ , in order to maximize  $\frac{1}{T} \sum_t r_t^\dagger$  while ensuring  $\mathbf{v}_{\text{avg}}^\dagger \in S$ . (We will use  $r_{\text{avg}}^\dagger$  to denote  $\frac{1}{T} \sum_t r_t^\dagger$ .) The proof for this special case will illustrate the main ideas required for proving regret bounds for the OSCP problem, i.e., the more general problem with ‘objective plus constraints’, over and above the techniques used in the previous section for the case of ‘only constraints’.

For this special case, Algorithm 2 reduces to the following:

---

**Algorithm 3** Linear objectives

---

Initialize  $\boldsymbol{\theta}_1$ .

**for all**  $t = 1, \dots, T$  **do**

    Choose option

$$(r_t^\dagger, \mathbf{v}_t^\dagger) = \arg \max_{(r, \mathbf{v}) \in A_t} r - 2Z\boldsymbol{\theta}_t \cdot \mathbf{v}.$$

    Choose  $\boldsymbol{\theta}_{t+1}$  by doing OCO update with  $g_t(\boldsymbol{\theta}) = \boldsymbol{\theta} \cdot \mathbf{v}_t^\dagger - h_S(\boldsymbol{\theta})$ , and domain  $W = \{\|\boldsymbol{\theta}\|_* \leq 1\}$ .

**end for**

---

**THEOREM 4.** *Given  $Z$  that satisfies Assumption 1, Algorithm 3 achieves the following regret bounds for OSCP with linear objective, in RP model:*

$$\begin{aligned} \mathbb{E}[\text{avg-regret}_1(T)] &\leq \frac{Z}{T} \cdot O(\mathcal{R}(T) + \mathcal{Q}(T)) \text{ and} \\ \mathbb{E}[\text{avg-regret}_2(T)] &\leq \frac{1}{T} \cdot O(\mathcal{R}(T) + \mathcal{Q}(T)). \end{aligned}$$

Here,  $\mathcal{Q}(T) = O(\|\mathbf{1}_d\| \sqrt{sT \log(d)})$ ,  $s = \max_{\mathbf{v} \in S} \max_j v_j$ , and  $\mathcal{R}(T)$  denotes the OCO regret for  $g_t(\cdot)$  over domain  $W$ .

*Proof.* Denote by  $(r_t^*, \mathbf{v}_t^*)$  the choice made by the offline optimal solution to satisfy request  $A_t$ . Then,

$$\mathbb{E}[r_t^*] = \text{OPT}, \text{ and } \mathbb{E}[\mathbf{v}_t^*] \in S,$$

where expectation is over  $A_t$  drawn uniformly at random from  $X_1, \dots, X_T$ .

Lemma 8 upper bounds  $\sum_t \mathbb{E}[2Zg_t(\boldsymbol{\theta}_t) - r_t^\dagger + r_t^*]$  by  $2Z\mathcal{Q}(T) = 2ZO(\|\mathbf{1}_d\| \sqrt{s \log(d)T})$ , using exactly the same line of argument as the proof of Lemma 6. Therefore, using  $\mathbb{E}[r_t^*] = \text{OPT}$ , the expected average reward obtained by the algorithm can be lower bounded as

$$\mathbb{E}[r_{\text{avg}}^\dagger] \geq \text{OPT} + \frac{2Z}{T} \sum_t \mathbb{E}[g_t(\boldsymbol{\theta}_t)] - \frac{2Z}{T} \mathcal{Q}(T).$$

As in the proof of Theorem 3, using Fenchel duality and OCO guarantees, it follows that  $d(\mathbf{v}_{\text{avg}}^\dagger, S) \leq \frac{1}{T} \sum_t g_t(\boldsymbol{\theta}_t) + \frac{1}{T} \mathcal{R}(T)$ , which gives,

$$\mathbb{E}[r_{\text{avg}}^\dagger] \geq \text{OPT} + (2Z)\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] - \frac{2Z}{T} \mathcal{R}(T) - \frac{2Z}{T} \mathcal{Q}(T). \quad (7)$$

Now, we use Assumption 1 to upper bound the reward obtained by the algorithm in terms of OPT and distance from set  $S$ . In particular, for  $\delta := \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]$ , since  $d(\mathbb{E}[\mathbf{v}_{\text{avg}}^\dagger], S) \leq \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] = \delta$ ,

$$\mathbb{E}[r_{\text{avg}}^\dagger] \leq \text{OPT}^\delta \leq \text{OPT} + Z\delta = \text{OPT} + Z \cdot \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]. \quad (8)$$

Combining inequalities (7) and (8), we obtain

$$\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \leq \frac{2}{T}\mathcal{R}(T) + \frac{2}{T}\mathcal{Q}(T),$$

and from (7), using the fact that  $\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \geq 0$ , we get that

$$\mathbb{E}[r_{\text{avg}}^\dagger] \geq \text{OPT} - \frac{2Z}{T} \cdot (\mathcal{R}(T) + \mathcal{Q}(T)).$$

This gives the theorem statement. □

LEMMA 8.  $\mathbb{E}[\sum_t 2Zg_t(\theta_t) - r_t^\dagger + r_t^*] \leq O(Z\|\mathbf{1}_d\|\sqrt{sT\log(d)})$ .

The proof of the above lemma follows exactly the same line of argument as the proof of Lemma 6. We omit it for brevity.

## 6. Online stochastic packing

Recall that the *online stochastic packing* problem is a special case of the online stochastic CP with linear objectives, with  $S = \{\mathbf{y} : \mathbf{y} \leq \frac{B}{T}\mathbf{1}\}$ . However, the performance of an algorithm for online stochastic packing is typically measured by *competitive ratio*, which is the ratio of total expected reward obtained by the online algorithm to the optimal solution or benchmark. The benchmarks in online packing are defined as sum of rewards, where as we defined OPT as the average reward. Therefore, in our notation, the competitive ratio for the online packing problem is given by  $\frac{\mathbb{E}[\sum_t r_t^\dagger]}{T\text{OPT}} = \frac{\mathbb{E}[\frac{1}{T}\sum_t r_t^\dagger]}{\text{OPT}}$ . The competitive ratio we obtain is  $1 - O(\epsilon)$ , for any  $\epsilon > 0$  such that  $\min\{B, T\text{OPT}\} \geq \log(d)/\epsilon^2$ .

Another important difference is that for online packing the budget is not allowed to be violated at all, while online CP allows a small violation of the constraint. A simple fix to make sure that budgets are not violated is to simply stop whenever a budget constraint is breached.<sup>2</sup> Another change we make to the algorithm is that we use a slightly different function in the OCO algorithm. We will use

$$g_t(\boldsymbol{\theta}) = (\mathbf{v}_t^\dagger - \frac{B}{T}\mathbf{1}) \cdot \boldsymbol{\theta}$$

over the domain  $\|\boldsymbol{\theta}\|_1 \leq 1, \boldsymbol{\theta} \geq \mathbf{0}$ . We use the Multiplicative Weight (MW) update algorithm as our OCO algorithm, which provides strong multiplicative guarantees (refer to Lemma 4 and Corollary 2).

Finally, as with the previous algorithms, we state the algorithm assuming we are given the parameter  $Z$ . We then show how to estimate  $Z$  to desired accuracy using only an  $O(\epsilon^2 \log(1/\epsilon))$  fraction of samples and solving an LP *only once* (in Lemma 12), assuming that  $\min\{B, \text{TOPT}\} \geq \frac{\log(d)}{\epsilon^2}$ .

We now state the algorithm below (Algorithm 4) for the online stochastic packing problem:

---

**Algorithm 4** Online Packing

---

Initialize  $\boldsymbol{\theta}_1 = \frac{1}{d+1} \mathbf{1}$ ,  $\mathbf{w}_1 = \mathbf{1}$ .

Initialize  $Z$  such that  $\frac{\text{TOPT}}{B} \leq Z \leq O(1) \frac{\text{TOPT}}{B}$ .

**for all**  $t = 1, \dots, T$  **do**

$$(r_t^\dagger, \mathbf{v}_t^\dagger) = \arg \max_{(r, \mathbf{v}) \in A_t} \{r - Z \boldsymbol{\theta}_t \cdot \mathbf{v}\}.$$

If, for some  $j = 1, \dots, d$ ,  $\sum_{t' \leq t} \mathbf{v}_{t'}^\dagger \cdot \mathbf{e}_j \geq B$  then EXIT.

Update  $\boldsymbol{\theta}_{t+1}$  using multiplicative weight update:

$$\forall j = 1..d, w_{t,j} = w_{t-1,j} (1 + \epsilon)^{\mathbf{v}_t^\dagger \cdot \mathbf{e}_j - B/T}$$

and

$$\forall j = 1..d, \boldsymbol{\theta}_{t+1,j} = \frac{w_{t,j}}{1 + \sum_{j'=1}^d w_{t,j'}}$$

**end for**

---

Strictly speaking, if we use the first few requests as samples to estimate  $Z$ , then we need to ignore these requests, and bound the error due to this. However, since the number of samples required is only  $O(\epsilon^2 \log(1/\epsilon))$  fraction of all requests, this error is quite small relative to the guarantee we obtain, which is a competitive ratio of  $1 - O(\epsilon)$ . We therefore ignore this error for the ease of presentation.

Let  $\tau$  be the stopping time of the algorithm. Denote by  $(r_t^*, \mathbf{v}_t^*)$  the choice made by the offline optimal solution to satisfy request  $A_t$ . We begin with the following lemma which is similar to Lemma 6.

LEMMA 9.

$$\sum_{t=1}^{\tau} \mathbb{E}[r_t^\dagger | \mathcal{F}_{t-1}] \geq \tau \text{OPT} + Z \sum_{t=1}^{\tau} \boldsymbol{\theta}_t \cdot \mathbb{E}[\mathbf{v}_t^\dagger - \mathbf{1} \frac{B}{T} | \mathcal{F}_{t-1}] - \sum_{t=1}^{\tau} Q(t),$$

where  $Q(t) = Z \|\mathbb{E}[\mathbf{v}_t^*] - \mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}]\| + \|\mathbb{E}[r_t^*] - \mathbb{E}[r_t^* | \mathcal{F}_{t-1}]\|$ .

*Proof.* If  $A_t$  is drawn uniformly at random from  $X_1, \dots, X_T$ , then  $\mathbb{E}[r_t^*] = \text{OPT}$ , and  $\mathbb{E}[\mathbf{v}_t^*] \leq \frac{B}{T} \mathbf{1}$ .

The algorithm chooses  $(r_t^\dagger, \mathbf{v}_t^\dagger) = \arg \max_{(r, \mathbf{v}) \in A_t} r - Z(\boldsymbol{\theta}_t \cdot \mathbf{v})$ . By the choice made by the algorithm

$$r_t^\dagger - Z(\boldsymbol{\theta}_t \cdot \mathbf{v}_t^\dagger) \geq r_t^* - Z(\boldsymbol{\theta}_t \cdot \mathbf{v}_t^*),$$

$$\begin{aligned} \mathbb{E}[r_t^\dagger - Z(\boldsymbol{\theta}_t \cdot \mathbf{v}_t^\dagger) | \mathcal{F}_{t-1}] &\geq \mathbb{E}[r_t^* | \mathcal{F}_{t-1}] - Z(\boldsymbol{\theta}_t \cdot \mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}]) \\ &\geq \mathbb{E}[r_t^*] - Z(\boldsymbol{\theta}_t \cdot \mathbb{E}[\mathbf{v}_t^*]) - Q(t) \\ &\geq \text{OPT} - Z\boldsymbol{\theta}_t \cdot \frac{B\mathbf{1}}{T} - Q(t). \end{aligned}$$

Summing above inequality for  $t = 1$  to  $\tau$  gives the lemma statement.  $\square$

LEMMA 10.

$$\sum_{t=1}^{\tau} \boldsymbol{\theta}_t \cdot (\mathbf{v}_t^\dagger - \frac{B}{T} \mathbf{1}) \geq B(1 - \epsilon - (1 + \epsilon) \frac{\tau}{T})^+ - \frac{\log(d+1)}{\epsilon},$$

where  $(a)^+$  denotes  $\max\{a, 0\}$ .

*Proof.* Recall that  $g_t(\boldsymbol{\theta}_t) = \boldsymbol{\theta}_t \cdot (\mathbf{v}_t^\dagger - \frac{B}{T} \mathbf{1})$ , therefore the LHS in the required inequality is  $\sum_{t=1}^{\tau} g_t(\boldsymbol{\theta}_t)$ .

We apply the regret bounds for the multiplicative weight update algorithm given by Corollary 2.

Now either  $\tau < T \frac{(1-\epsilon)}{(1+\epsilon)}$ , which means that the algorithm aborted due to constraint violation, i.e.,  $\sum_{t=1}^{\tau} (\mathbf{v}_t^\dagger \cdot \mathbf{e}_j) \geq B$  for some  $j$  at the stopping time  $\tau$ . And, we have that  $\sum_{t=1}^{\tau} g_t(\mathbf{e}_j) \geq B - \frac{\tau B}{T}$ . Corollary 2 (with  $M = 1, R = B/T$ ) guarantees that

$$\begin{aligned} \sum_{t=1}^{\tau} g_t(\boldsymbol{\theta}_t) &\geq (1 - \epsilon) \left( \sum_{\geq 0} g_t(\mathbf{e}_j) \right) + (1 + \epsilon) \left( \sum_{< 0} g_t(\mathbf{e}_j) \right) - \frac{\log(d+1)}{\epsilon} \\ &\geq (1 - \epsilon) \left( \sum_t g_t(\mathbf{e}_j) \right) + 2\epsilon \left( \sum_{< 0} g_t(\mathbf{e}_j) \right) - \frac{\log(d+1)}{\epsilon} \\ &\geq (1 - \epsilon) \left( B - \frac{\tau B}{T} \right) - 2\epsilon \frac{\tau B}{T} - \frac{\log(d+1)}{\epsilon} \\ &= (1 - \epsilon)B - (1 + \epsilon) \frac{\tau B}{T} - \frac{\log(d+1)}{\epsilon} \end{aligned}$$

Or,  $\tau \geq T \frac{(1-\epsilon)}{(1+\epsilon)}$ , in which case  $(1 - \epsilon - \frac{\tau(1+\epsilon)}{T})^+ = 0$ , and the right hand side in the lemma statement is equal to  $-\frac{\log(d+1)}{\epsilon}$ . The inequality is trivially obtained from Corollary 2, which guarantees,

$$\sum_{t=1}^{\tau} g_t(\boldsymbol{\theta}_t) \geq -\frac{\log(d+1)}{\epsilon}$$

$\square$

Now, we are ready to prove Theorem 2, which states that Algorithm 4 achieves a competitive ratio of  $1 - O(\epsilon)$ , given  $\min\{B, T\text{OPT}\} \geq \frac{\log(d)}{\epsilon^2}$  for the online stochastic packing problem in RP model.

**Proof of Theorem 2.** Substituting the inequality from Lemma 10 in Lemma 9, we get

$$\sum_{t=1}^{\tau} \mathbb{E}[r_t^\dagger | \mathcal{F}_{t-1}] \geq \tau \text{OPT} + ZB \left( (1 - \epsilon - (1 + \epsilon) \frac{\tau}{T})^+ \right) - Z \frac{\log(d+1)}{\epsilon} - \sum_{t=1}^{\tau} Q(t)$$

Now, using  $Z \leq O(1) \frac{T\text{OPT}}{B}$  and  $B \geq \frac{\log(d)}{\epsilon^2}$ , we get

$$Z \frac{\log(d+1)}{\epsilon} \leq O(1) \frac{T\text{OPT}}{B} \frac{\log(d+1)}{\epsilon} = O(\epsilon) T\text{OPT}.$$

Also,  $Z \geq \frac{T\text{OPT}}{B}$ . Substituting in above,

$$\sum_{t=1}^{\tau} \mathbb{E}[r_t^\dagger | \mathcal{F}_{t-1}] \geq \tau \text{OPT} + \left( (1 - \epsilon - (1 + \epsilon) \frac{\tau}{T})^+ \right) T\text{OPT} - O(\epsilon) T\text{OPT} - \sum_{t=1}^{\tau} Q(t)$$

Now, either  $\tau \geq T \frac{(1-\epsilon)}{(1+\epsilon)}$ , in which case above gives

$$\begin{aligned} \sum_{t=1}^{\tau} \mathbb{E}[r_t^\dagger | \mathcal{F}_{t-1}] &\geq \frac{(1-\epsilon)}{(1+\epsilon)} T\text{OPT} - O(\epsilon) T\text{OPT} - \sum_{t=1}^{\tau} Q(t) \\ &= (1 - O(\epsilon)) T\text{OPT} - O(\epsilon) T\text{OPT} - \sum_{t=1}^{\tau} Q(t) \end{aligned}$$

Otherwise,  $(1 - \epsilon - (1 + \epsilon) \frac{\tau}{T})^+ = (1 - \epsilon - (1 + \epsilon) \frac{\tau}{T})$ , and from above we get,

$$\begin{aligned} \sum_{t=1}^{\tau} \mathbb{E}[r_t^\dagger | \mathcal{F}_{t-1}] &\geq \tau \text{OPT} + \left( 1 - \epsilon - (1 + \epsilon) \frac{\tau}{T} \right) T\text{OPT} - O(\epsilon) T\text{OPT} - \sum_{t=1}^{\tau} Q(t) \\ &= (1 - \epsilon) T\text{OPT} - \epsilon \tau \text{OPT} - O(\epsilon) T\text{OPT} - \sum_{t=1}^{\tau} Q(t) \\ &= (1 - \epsilon) T\text{OPT} - O(\epsilon) T\text{OPT} - \sum_{t=1}^{\tau} Q(t). \end{aligned}$$

Then, taking expectation on both sides,  $\mathbb{E}[\sum_{t=1}^{\tau} r_t^\dagger] \geq (1 - O(\epsilon)) T\text{OPT} - \mathbb{E}[\sum_{t=1}^{\tau} Q(t)]$ .

Just like in the proof of Lemma 6, we can bound  $\mathbb{E}[\sum_{t=1}^{\tau} Q(t)] \leq Z \|\mathbf{1}_{d+1}\|_{\infty} \sqrt{sT \log(d+1)}$  which is  $O(\epsilon) T\text{OPT}$ , using the fact that for  $S = \{\mathbf{y} : \mathbf{y} \leq \frac{B}{T} \mathbf{1}\}$ , the parameter  $s = \max_{j, \mathbf{y} \in S} y_j = \frac{B}{T}$ ,  $\|\mathbf{1}_{d+1}\|_{\infty} = 1$ , and that  $Z \leq O(1) \frac{T\text{OPT}}{B}$ ,  $\epsilon \geq \sqrt{\frac{\log(d)}{B}}$ . This completes the proof.  $\square$

We now show how to compute a  $Z$  as required using the first  $O(\epsilon^2 \log(1/\epsilon))$  requests as samples. For convenience, let  $\text{OPT}_{\text{SUM}} := T\text{OPT}$  denote the optimum for the sum. We first state a lemma that relates the optimum value of an offline packing instance to the optimum value on a sample of the requests. The proof of this is along the lines of a similar lemma (Lemma 14) in Devanur et al. (2011b), and we present the proof in Appendix G for the sake of completeness.

LEMMA 11. For all  $\rho \in (0, 1]$ , there exists  $\eta = O\left(\sqrt{\log\left(\frac{d}{\rho}\right)}\right)$  such that for all  $\delta \in (0, 1]$ , given a random sample of  $\delta T$  requests, one can compute a quantity  $\hat{OPT}$  such that with probability  $1 - \rho$ ,

1.  $\hat{OPT} \geq OPT_{\text{SUM}} - \eta\sqrt{OPT_{\text{SUM}}/\delta}$ .
2.  $\frac{\hat{OPT}}{1+\eta/\sqrt{\delta B}} \leq OPT_{\text{SUM}} + \eta\sqrt{OPT_{\text{SUM}}/\delta}$ .

LEMMA 12. Given a random sample of  $O(\epsilon^2 \log(1/\epsilon))$  fraction of requests, one can compute a quantity  $Z$  such that with probability at least  $1 - \epsilon^2$ ,

$$\frac{OPT_{\text{SUM}}}{B} \leq Z \leq \frac{9}{2} \frac{OPT_{\text{SUM}}}{B}.$$

*Proof.* We use Lemma 11 with  $\rho = \epsilon^2$  and  $\delta = 4\eta^2\epsilon^2/\log(d)$ . Then, from the assumption that  $\min\{B, OPT_{\text{SUM}}\} \geq \log(d)/\epsilon^2$ , we have that  $\delta \geq 4\eta^2/OPT_{\text{SUM}}$ , and  $\delta \geq 4\eta^2/B$ . Therefore, we get that with probability at least  $1 - \epsilon^2$ ,

$$\begin{aligned} \hat{OPT} &\geq OPT_{\text{SUM}} - \eta\sqrt{OPT_{\text{SUM}}/\delta} \\ &\geq OPT_{\text{SUM}} - OPT_{\text{SUM}}/2 = OPT_{\text{SUM}}/2. \end{aligned}$$

Also,

$$\begin{aligned} \hat{OPT} &\leq (1 + \eta/\sqrt{\delta B})(OPT_{\text{SUM}} + \eta\sqrt{OPT_{\text{SUM}}/\delta}) \\ &\leq \frac{3}{2}(OPT_{\text{SUM}} + \frac{1}{2}OPT_{\text{SUM}}) \\ &\leq \frac{9}{4}OPT_{\text{SUM}}. \end{aligned}$$

Therefore  $Z := 2\hat{OPT}/B$  satisfies the conclusion of the lemma. Finally, note that  $\delta = 4\eta^2\epsilon^2/\log(d) = O(\epsilon^2 \log(d)/\log(d)) = O(\epsilon^2 \log(1/\epsilon))$ .  $\square$

## 7. Stronger bounds for smooth functions

We show that when  $f$  is a strongly smooth function, and, instead of distance function, a strongly smooth function is used to measure regret in constraint violation, then stronger regret bounds of  $O(\frac{\log T}{T})$  can be achieved in IID case.

More precisely, consider the following *smooth* version of OSCP.



DEFINITION 4. [ONLINE STOCHASTIC SMOOTH CONVEX PROGRAMMING] Let  $f$  be a  $\beta$ -smooth concave function. And, let  $h$  be a  $\beta$ -smooth convex function. At time  $t$ , the algorithm needs to choose  $\mathbf{v}_t^\dagger \in A_t$  to minimize regret defined as

$$\text{avg-regret}_1(T) := f(\mathbf{v}_{\text{avg}}^*) - f(\mathbf{v}_{\text{avg}}^\dagger),$$

$$\text{avg-regret}_2(T) := h(\mathbf{v}_{\text{avg}}^\dagger).$$

Here,  $\mathbf{v}_{\text{avg}}^* = \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^*$ ,  $\mathbf{v}_{\text{avg}}^\dagger = \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^\dagger$ . Also, assume that there exist  $\mathbf{v}_t \in A_t$  for all  $t$ , such that  $h(\frac{1}{T} \sum_t \mathbf{v}_t) = 0$ .

Note that we do not require Lipschitz condition for  $f$  or  $h$ .

We provide two algorithms for this problem. The first algorithm is a simple extension of Algorithm 2. We show that, under certain technical conditions, this algorithm achieves stronger regret bounds of  $\tilde{O}(\frac{\log T}{T})$  in smooth case. Intuitively, this is because as discussed in Section 3, the dual of strongly smooth functions is strongly convex, and for strongly convex/concave functions, stronger logarithmic regret guarantees are provided by online learning algorithms. The second algorithm is a novel primal algorithm based on Frank-Wolfe algorithm for solving smooth convex programs.

### 7.1. Primal-dual algorithm based on Algorithm 2

We make an additional assumption.

ASSUMPTION 2. Let  $\nabla_f$  and  $\nabla_g$  denote the set of gradients of functions  $f$  and  $g$ , respectively, on domain  $[0, 1]^d$ , i.e.,

$$\nabla_f = \{\nabla f(\mathbf{x}) : \mathbf{x} \in [0, 1]^d\}, \text{ and ,}$$

$$\nabla_g = \{\nabla g(\mathbf{x}) : \mathbf{x} \in [0, 1]^d\}.$$

Assume that the sets  $cl(\nabla_f)$  and  $cl(\nabla_g)$  are convex and easy to project upon. Here  $cl(S)$  denotes the closure of set  $S$ .

This assumption is true for many natural concave utility and convex risk functions, in particular, for all separable smooth functions. Now, an algorithm similar to Algorithm 2 can be used for this

problem. One change we make is that we perform online learning for  $g_t$  and  $\psi_t$  on domain  $\nabla_g$  and  $\nabla_f$ , respectively, which is possible because from Assumption 2, these domains are convex and easy to project upon.

---

**Algorithm 5** Algorithm for smooth case based on Algorithm 2
 

---

Initialize  $\theta_1, \phi_1$ .

**for all**  $t = 1, \dots, T$  **do**

  Choose vector

$$\mathbf{v}_t^\dagger = \arg \max_{\mathbf{v} \in A_t} -\phi_t \cdot \mathbf{v} - 2Z\theta_t \cdot \mathbf{v}.$$

  Choose  $\theta_{t+1}$  by doing an OCO update for  $g_t(\theta) = \theta \cdot \mathbf{v}_t^\dagger - h_S(\theta)$  over domain  $\nabla_g$ .

  Choose  $\phi_{t+1}$  by doing an OCO update for  $\psi_t(\phi) = \phi \cdot \mathbf{v}_t^\dagger - (-f)^*(\phi)$  over domain  $\nabla_f$ .

**end for**

---

**THEOREM 5.** *Under Assumption 2, and given  $Z$  that satisfies Assumption 1, Algorithm 5 achieves the following regret for the Online Smooth Convex Programming problem, in the stochastic IID input model.*

$$\begin{aligned} \mathbb{E}[\text{avg-regret}_1(T)] &= Z \cdot O\left(\frac{C \log(T)}{T}\right), \\ \mathbb{E}[\text{avg-regret}_2(T)] &= O\left(\frac{C \log(T)}{T}\right), \end{aligned}$$

where  $C = \beta \|\mathbf{1}_d\|^2$ .

*Proof.* The proof follows from the proof of Theorem 1 on observing that stronger OCO regret bounds of  $O(\log(T))$  are available for strongly convex functions. More precisely, in case of IID inputs, the proof of Theorem 1 can be followed as it is to achieve the following regret bounds. (These are same as in the detailed statement of Theorem 1, provided in Appendix E, but with  $\mathcal{Q}(T) = 0$  due to IID assumption.)

$$\begin{aligned} \mathbb{E}[\text{avg-regret}_1(T)] &\leq \frac{Z}{T} \cdot O(\mathcal{R}(T)) + O\left(\frac{\mathcal{R}'(T)}{T}\right), \\ \mathbb{E}[\text{avg-regret}_2(T)] &\leq \frac{1}{T} \cdot O(\mathcal{R}(T)) + \frac{1}{Z} O\left(\frac{\mathcal{R}'(T)}{T}\right). \end{aligned}$$

Here  $\mathcal{R}(T)$  is OCO regret for the problem of maximizing concave function  $g_t(\theta) = \theta \cdot \mathbf{v}_t - h^*(\theta)$ ,  $\mathcal{R}'(T)$  is OCO regret for the problem of maximizing concave function  $\psi_t(\phi) = \phi \cdot \mathbf{v}_t - (-f)^*(\phi)$ . Now, using

Lemma 2, given that  $h$  and  $f$  are  $\beta$ -strongly smooth,  $g_t$  and  $\psi_t$  are  $\frac{1}{\beta}$ -strongly concave over domain  $\nabla_g$  and  $\nabla_f$  respectively. Also, the gradient of these functions is some  $\mathbf{v} \in [0, 1]^d$ , so that the norms of gradients are bounded by  $\|\mathbf{1}_d\|$ .

Therefore, using online learning guarantees for smooth functions from Lemma 5, along with  $G = \|\mathbf{1}_d\|, H = 1/\beta$ , we get  $\mathcal{R}(T) = O(\|\mathbf{1}_d\|^2 \beta \log T)$ , and  $\mathcal{R}'(T) = O(\|\mathbf{1}_d\|^2 \beta \log T)$ . The theorem statement is obtained by substituting these OCO regret bounds in above.  $\square$

In above, observe that Assumption 2 was required because Lemma 2 provided strong convexity of  $g_t(\cdot)$  and  $\psi_t(\cdot)$  only on the domains  $\nabla_g$  and  $\nabla_f$ , respectively. We conjecture that it is possible to remove this assumption to get similar regret guarantees for the smooth case.

## 7.2. Primal algorithm based on Frank-Wolfe algorithm

For simplicity let us first consider the problem with only objective function and no constraints. The overall goal of online algorithm is to choose  $\mathbf{v}_t^\dagger \in A_t$  at time  $t$  in order to maximize  $f(\mathbf{v}_{\text{avg}}^\dagger)$ , where  $\mathbf{v}_{\text{avg}}^\dagger = \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^\dagger$ , and  $f$  is  $\beta$ -smooth concave function. We compare to the optimal solution of the (offline) expected value instance where the optimal solution solves

$$\max_{\{\mathbf{v}_A \in A\}} \mathbb{E}_A[f(\mathbf{v}_A)]$$

Let  $\mathbf{v}_t^* \in A_t$  denotes the vector chosen from set  $A_t$  by the optimal solution, and let  $\mathbf{v}_{\text{avg}}^* = \mathbb{E}[\mathbf{v}_t^*]$ , where expectation is over randomly (i.i.d.) generated set  $A_t$ . Regret is defined as

$$\text{avg-regret}_1(T) = f\left(\frac{1}{T} \sum_t \mathbf{v}_t^*\right) - f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^\dagger\right).$$

And, expected regret

$$\mathbb{E}[\text{avg-regret}_1(T)] = \mathbb{E}\left[f\left(\frac{1}{T} \sum_t \mathbf{v}_t^*\right) - f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t^\dagger\right)\right] \leq f(\mathbf{v}_{\text{avg}}^*) - \mathbb{E}[f(\mathbf{v}_{\text{avg}}^\dagger)].$$

For this special case, we obtain the following primal algorithm based on Frank-Wolfe algorithm.

---

**Algorithm 6** Algorithm for smooth case based on Frank-Wolfe (no constraints)
 

---

**for all**  $t = 1, \dots, T$  **do**

 Let  $\mathbf{x}_t = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbf{v}_\tau^\dagger$ . Use arbitrary  $\mathbf{x}_1$ .

Choose vector

$$\mathbf{v}_t^\dagger = \arg \max_{\mathbf{v} \in A_t} \nabla f(\mathbf{x}_t) \cdot \mathbf{v}$$

**end for**


---

**THEOREM 6.** *Algorithm 6 achieves the following regret bound for Online Smooth Convex Programming in the IID model of stochastic inputs:*

$$\mathbb{E}[\text{avg-regret}_1(T)] = O\left(\frac{\beta \log(T)}{T}\right).$$

*Proof.* Let  $\Delta_t := f(\mathbf{v}_{\text{avg}}^*) - \mathbb{E}[f(\mathbf{x}_t)]$ , where expectation is over random (i.i.d.) generation of sets  $A_t$ .

We prove that  $\mathbb{E}[\Delta_T] \leq \frac{\beta \log(2T)}{2T}$ . (The base of the log is 2.) This will imply the required expected regret bound. By concavity,

$$f(\mathbf{v}_{\text{avg}}^*) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) \cdot (\mathbf{v}_{\text{avg}}^* - \mathbf{x}_t).$$

Now,  $\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] = \mathbf{v}_{\text{avg}}^*$ , where expectation is over randomly (i.i.d.) generated set  $A_t$ , given any history  $\mathcal{F}_{t-1}$ . Also, by definition,  $\mathbf{x}_t$  is fixed by the history  $\mathcal{F}_{t-1}$ . Therefore, above implies

$$f(\mathbf{v}_{\text{avg}}^*) \leq \mathbb{E}[(f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) \cdot (\mathbf{v}_t^* - \mathbf{x}_t)) | \mathcal{F}_{t-1}]$$

Now, since  $\mathbf{v}_t^\dagger$  was chosen from set  $A_t$  to maximize  $\nabla f(\mathbf{x}_t) \cdot \mathbf{v}_t^\dagger$ , we have that

$$\begin{aligned} f(\mathbf{v}_{\text{avg}}^*) &\leq \mathbb{E}[f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) \cdot (\mathbf{v}_t^\dagger - \mathbf{x}_t) | \mathcal{F}_{t-1}] \\ &= f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t) \cdot (\mathbb{E}[\mathbf{v}_t^\dagger | \mathcal{F}_{t-1}] - \mathbf{x}_t). \end{aligned} \tag{9}$$

Now,  $\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{1}{(t+1)}(\mathbf{v}_t^\dagger - \mathbf{x}_t)$ . Using  $\beta$ -smoothness of  $f(\cdot)$ , we have that

$$f(\mathbf{x}_{t+1}) \geq f(\mathbf{x}_t) + \frac{1}{(t+1)} \nabla f(\mathbf{x}_t) \cdot (\mathbf{v}_t^\dagger - \mathbf{x}_t) - \frac{\beta}{2(t+1)^2}.$$

Taking conditional expectation on both sides, and substituting from Equation (9),

$$\begin{aligned} \mathbb{E}[f(\mathbf{x}_{t+1})|\mathcal{F}_{t-1}] &\geq f(\mathbf{x}_t) + \frac{1}{(t+1)} \nabla f(\mathbf{x}_t) \cdot (\mathbb{E}[\mathbf{v}_t^\dagger|\mathcal{F}_{t-1}] - \mathbf{x}_t) - \frac{\beta}{2(t+1)^2} \\ &\geq f(\mathbf{x}_t) + \frac{1}{(t+1)} (f(\mathbf{v}_{\text{avg}}^*) - f(\mathbf{x}_t)) - \frac{\beta}{2(t+1)^2}. \end{aligned}$$

Subtracting left and right hand side from  $f(\mathbf{v}_{\text{avg}}^*)$ , and taking expectations, we get

$$\begin{aligned} \Delta_{t+1} &\leq \Delta_t - \frac{1}{(t+1)} \Delta_t + \frac{\beta}{2(t+1)^2} \\ &= \frac{t}{(t+1)} \Delta_t + \frac{\beta}{2(t+1)^2}. \end{aligned} \tag{10}$$

Recall that we wish to show that  $\Delta_t \leq \beta \log(2t)/2t$ . The rest of the proof is by induction on  $t$ . For the base case, we note that we can still use (10) with  $t=0$  and an arbitrary  $\mathbf{x}_1$  which is used to choose  $\mathbf{v}_1^\dagger$ . This gives us that  $\Delta_1 \leq \beta/2$ . The inductive step for  $t+1$  follows from (10) and the inductive hypothesis for  $t$  if

$$\begin{aligned} \frac{t}{(t+1)} \cdot \frac{\beta \log(2t)}{2t} + \frac{\beta}{2(t+1)^2} &\leq \frac{\beta \log(2(t+1))}{2(t+1)} \\ \Leftrightarrow \log(t) + \frac{1}{t+1} &\leq \log(t+1) \\ \Leftrightarrow \frac{1}{t+1} &\leq \log(1 + \frac{1}{t}). \end{aligned}$$

The last inequality follows from the fact that for any  $a > 0$ ,  $\log(1+a) > \frac{a}{1+a}$ , by setting  $a = 1/t$ . This completes the induction. Therefore,  $\Delta_T = f(\mathbf{v}_{\text{avg}}^*) - f(\mathbf{x}_T) = f(\mathbf{v}_{\text{avg}}^*) - f(\mathbf{v}_{\text{avg}}^\dagger) \leq \frac{\beta \log(2T)}{2T}$ .  $\square$

For the Smooth online stochastic CP problem, the algorithm needs to be modified using the parameter  $Z$  to combine constraints and objective.

---

**Algorithm 7** Frank-Wolfe based algorithm for smooth functions

---

**for all**  $t = 1, \dots, T$  **do**

Let  $\mathbf{x}_t = \frac{1}{t-1} \sum_{\tau=1}^{t-1} \mathbf{v}_\tau^\dagger$ . Use arbitrary  $\mathbf{x}_1$ .

Choose vector

$$\mathbf{v}_t^\dagger = \arg \max_{\mathbf{v} \in A_t} \nabla f(\mathbf{x}_t) \cdot \mathbf{v} - Z(\nabla h(\mathbf{x}_t) \cdot \mathbf{v})$$

**end for**

---

## Appendix A: Concentration Inequalities

LEMMA 13. *Hoeffding (1963), Theorem 4*

Let  $\mathcal{X} = (x_1, \dots, x_N)$  be a finite population of  $N$  real points,  $X_1, \dots, X_n$  denote a random sample without replacement from  $\mathcal{X}$ , and  $Y_1, \dots, Y_n$  denote a random sample with replacement from  $\mathcal{X}$ . If  $\ell: \mathbb{R} \rightarrow \mathbb{R}$  is continuous and convex, then

$$\mathbb{E}[\ell(\sum_{t=1}^n X_t)] \leq \mathbb{E}[\ell(\sum_{t=1}^n Y_t)].$$

LEMMA 14. *Hoeffding (1963)* Let  $\mathcal{X} = (x_1, \dots, x_N)$  be a finite population of  $N$  real points,  $X_1, \dots, X_n$  denote a random sample without replacement from  $\mathcal{X}$ . Let  $a = \min_{1 \leq i \leq N} x_i$ ,  $b = \max_{1 \leq i \leq N} x_i$  and  $\mu = \frac{1}{N} \sum_{i=1}^N X_i$ . Then, for all  $\epsilon > 0$ ,

$$\Pr\left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \geq \epsilon\right) \leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right).$$

LEMMA 15. (*Multiplicative version*) Let  $\mathcal{X} = (x_1, \dots, x_N)$  be a finite population of  $N$  real points, and  $X_1, \dots, X_n$  denote a random sample without replacement from  $\mathcal{X}$ . Let  $a = \min_{1 \leq i \leq N} x_i$ ,  $b = \max_{1 \leq i \leq N} x_i$  and  $\mu = \frac{n}{N} \sum_{i=1}^N X_i$ . Then, for all  $\epsilon > 0$ ,

$$\Pr\left(\left|\sum_{i=1}^n X_i - \mu\right| \geq \epsilon\mu\right) \leq \exp\left(-\frac{\mu\epsilon^2}{3(b-a)^2}\right).$$

COROLLARY 3. (to Lemma 15) Let  $\mathcal{X} = (x_1, \dots, x_N)$  be a finite population of  $N$  real points, and  $X_1, \dots, X_n$  denote a random sample without replacement from  $\mathcal{X}$ . Let  $a = \min_{1 \leq i \leq N} x_i$ ,  $b = \max_{1 \leq i \leq N} x_i$  and  $\mu = \frac{n}{N} \sum_{i=1}^N X_i$ . Then, for all  $\rho > 0$ , with probability at least  $1 - \rho$ ,

$$\left|\sum_{i=1}^n X_i - \mu\right| \leq (b-a)\sqrt{3\mu \log(1/\rho)}$$

*Proof.* Given  $\rho > 0$ , use Lemma 15 with

$$\epsilon = (b-a)\sqrt{\frac{3\log(1/\rho)}{\mu}},$$

to get that the probability of the event  $|\sum_{i=1}^n X_i - \mu| > \epsilon\mu = (b-a)\sqrt{3\mu \log(1/\rho)}$  is at most

$$\exp\left(-\frac{\mu\epsilon^2}{3(b-a)^2}\right) = \exp(-\log(1/\rho)) = \rho.$$

□

LEMMA 16. *Kleinberg et al. (2008), Babaioff et al. (2012), Badanidiyuru et al. (2013)* Consider a probability distribution with values in  $[0, 1]$ , and expectation  $\nu$ . Let  $\hat{\nu}$  be the average of  $N$  independent samples from this distribution. Then, with probability at least  $1 - e^{-\Omega(\gamma)}$ , for all  $\gamma > 0$ ,

$$|\hat{\nu} - \nu| \leq \text{rad}(\hat{\nu}, N) \leq 3\text{rad}(\nu, N), \quad (11)$$

where  $\text{rad}(\nu, N) = \sqrt{\frac{\nu}{N}} + \frac{\nu}{N}$ . More generally this result holds if  $X_1, \dots, X_N \in [0, 1]$  are random variables,  $N\hat{\nu} = \sum_{t=1}^N X_t$ , and  $N\nu = \sum_{t=1}^N \mathbb{E}[X_t | X_1, \dots, X_{t-1}]$ .

## Appendix B: Preliminaries

### B.1. Strong smoothness/Strong convexity duality.

*Proof of Lemma 2* Given  $h$  is convex and  $\beta$ -strong smooth with respect to norm  $\|\cdot\|$ . We prove that  $h^*$ , defined as

$$h^*(\boldsymbol{\theta}) = \max_{\mathbf{y} \in [0,1]^d} \{\mathbf{y} \cdot \boldsymbol{\theta} - h(\mathbf{y})\},$$

is  $\frac{1}{\beta}$ -strongly convex with respect to norm  $\|\cdot\|_*$  on domain  $\nabla_h = \{\nabla h(\mathbf{x}) : \mathbf{x} \in [0,1]^d\}$ .

For any  $\boldsymbol{\theta}, \boldsymbol{\phi} \in \nabla_h$ ,  $\boldsymbol{\theta} = \nabla h(\mathbf{z}), \boldsymbol{\phi} = \nabla h(\mathbf{x})$  for some  $\mathbf{z}, \mathbf{x} \in [0,1]^d$ . And, therefore,

$$\begin{aligned} h^*(\boldsymbol{\theta}) - h^*(\boldsymbol{\phi}) - \mathbf{x} \cdot (\boldsymbol{\theta} - \boldsymbol{\phi}) &= h^*(\nabla h(\mathbf{z})) - h^*(\nabla h(\mathbf{x})) - \mathbf{x} \cdot (\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) \\ &= \mathbf{z} \cdot \nabla h(\mathbf{z}) - h(\mathbf{z}) - (\mathbf{x} \cdot \nabla h(\mathbf{x}) - h(\mathbf{x})) - \mathbf{x} \cdot (\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) \\ &= \mathbf{z} \cdot \nabla h(\mathbf{z}) - h(\mathbf{z}) + h(\mathbf{x}) - \mathbf{x} \cdot \nabla h(\mathbf{z}) \\ &= (\mathbf{z} - \mathbf{x}) \cdot (\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) - (h(\mathbf{z}) - h(\mathbf{x}) - \nabla h(\mathbf{x})(\mathbf{z} - \mathbf{x})) \\ &= (\mathbf{z} - \mathbf{x}) \cdot (\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) - g(\mathbf{z} - \mathbf{x}), \end{aligned} \tag{12}$$

where we define

$$g(\mathbf{y}) := h(\mathbf{x} + \mathbf{y}) - h(\mathbf{x}) - (\nabla h(\mathbf{x})) \cdot \mathbf{y}.$$

Now, for any  $\varphi$ ,

$$\begin{aligned} g^*(\varphi) &:= \sup_{\mathbf{y}} \varphi \cdot \mathbf{y} - g(\mathbf{y}) \\ &= \varphi \cdot \mathbf{y}^* - g(\mathbf{y}^*) \end{aligned}$$

where  $\mathbf{y}^*$  is such that  $\varphi = \nabla g(\mathbf{y}^*) = \nabla h(\mathbf{x} + \mathbf{y}^*) - \nabla h(\mathbf{x})$ . Therefore, for  $\varphi = \nabla h(\mathbf{z}) - \nabla h(\mathbf{x})$ ,  $\mathbf{y}^* = \mathbf{z} - \mathbf{x}$ , so that,

$$g^*(\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) = (\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) \cdot (\mathbf{z} - \mathbf{x}) - g(\mathbf{z} - \mathbf{x}).$$

Substituting in (12), we get

$$\begin{aligned} &h^*(\boldsymbol{\theta}) - h^*(\boldsymbol{\phi}) - \mathbf{x} \cdot (\boldsymbol{\theta} - \boldsymbol{\phi}) \\ &= g^*(\nabla h(\mathbf{z}) - \nabla h(\mathbf{x})) \text{ [Nikhil:] Why is this?} \\ &= g^*(\boldsymbol{\theta} - \boldsymbol{\phi}) \end{aligned}$$

By smoothness assumption,  $g(\mathbf{y}) \leq \frac{\beta}{2} \|\mathbf{y}\|^2$ . This implies that  $g^*(\boldsymbol{\theta}) \geq \frac{1}{2\beta} \|\boldsymbol{\theta}\|_*^2$  because the conjugate of  $\beta$  times half squared norm is  $1/\beta$  times half squared of the dual norm. This gives

$$h^*(\boldsymbol{\theta}) - h^*(\boldsymbol{\phi}) - \mathbf{x} \cdot (\boldsymbol{\theta} - \boldsymbol{\phi}) \geq \frac{1}{2\beta} \|\boldsymbol{\theta} - \boldsymbol{\phi}\|_*^2.$$

This completes the proof. □

## B.2. Online learning.

A popular algorithm for OCO is the online mirror descent (OMD) algorithm. The OMD algorithm with regularizer  $R(\boldsymbol{\theta})$  uses the following fast update rule to select player's decision  $\boldsymbol{\theta}_{t+1}$  for this problem:

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &= \arg \max_{\boldsymbol{\theta} \in W} \frac{1}{\eta} R(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{y}_{t+1}, \text{ where} \\ \mathbf{y}_{t+1} &= \mathbf{y}_t - z_t, \text{ and } z_t \in \partial g_t(\boldsymbol{\theta}_t)\end{aligned}\tag{13}$$

The maximization problem in above is particularly simple when domain  $W$  is of form  $\|\boldsymbol{\theta}\| \leq \gamma$ , and this is the main use case of this algorithm in this paper. Further, for domain  $W$  of form  $\|\boldsymbol{\theta}\|_2 \leq L$ , and  $R(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2$ , this simply becomes online gradient descent. OMD has the following guarantees for this problem:

LEMMA 17. *Shalev-Shwartz (2012)*

$$\mathcal{R}(T) \leq \frac{D}{\eta} + \eta T G^2,$$

where  $D = (\max_{\boldsymbol{\theta}''} R(\boldsymbol{\theta}'') - \min_{\boldsymbol{\theta}' \in W} R(\boldsymbol{\theta}'))$ ,  $\frac{1}{T} \sum_{t=1}^T \|z_t\|_2^2 \leq G$  for  $z_t \in \partial g_t(\boldsymbol{\theta}_t)$ , and  $R$  is a 1-strongly-convex function with respect to norm  $\|\cdot\|_*$ .

Now, to derive Corollary 1, observe that for  $W = \{\|\boldsymbol{\theta}\|_2 \leq L\}$ , Euclidean regularizer  $R(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2$  gives  $\mathcal{R}(T) \leq LG\sqrt{T}$ , with  $G^2 = d \geq \frac{1}{T} \sum_{t=1}^T \|z_t\|_2^2$ , when  $z_t \in [0, 1]^d$ . And, for  $W = \{\|\boldsymbol{\theta}\|_1 \leq L, \boldsymbol{\theta} > 0\}$ , entropic regularizer  $R(\boldsymbol{\theta}) = \sum_i \boldsymbol{\theta}_i \log \boldsymbol{\theta}_i$  gives  $\mathcal{R}(T) \leq G\sqrt{LT \log(d)}$ , where  $G^2 = 1 \geq \frac{1}{T} \sum_{t=1}^T \|z_t\|_\infty^2$ , when  $z_t \in [0, 1]^d$ .

## Appendix C: Sampling without replacement bounds for Section 4

**Proof of Equation (5).** Let  $\boldsymbol{\omega} = \mathbb{E}[\mathbf{w}_{t,\sigma}] = \mathbb{E}[\mathbf{v}_t^*]$ . To bound the quantity  $\mathbb{E}[\|\mathbf{w}_{t,\pi} - \boldsymbol{\omega}\|]$ , note that  $\mathbf{w}_{t,\pi}$  can be viewed as the average of  $t$  vectors sampled uniformly *without replacement* from the ground set  $\{\mathbf{v}_{X_1}, \dots, \mathbf{v}_{X_T}\}$  of  $T$  vectors.

Now, let  $w_{t,\pi,j}$  denote the  $j^{\text{th}}$  component of vector  $\mathbf{w}_{t,\pi}$ . Then, by applying concentration bounds from Corollary 3, we get that

$$|w_{t,\pi,j} - \omega_j| \leq \sqrt{\frac{3\omega_j \log(d/\rho)}{t}},$$

with probability  $1 - \frac{\rho}{d}$  for all  $\rho \in (0, 1)$ . From the condition  $\boldsymbol{\omega} = \mathbb{E}[\mathbf{v}_t^*] \in S$ , we have  $\omega_j \leq \max_{v \in S} v_j \leq s$ . Taking union bound over  $d$ , for every  $\rho \in (0, 1)$ , we have that with probability  $1 - \rho$ ,

$$\|\mathbf{w}_{t,\pi} - \boldsymbol{\omega}\| \leq \|\mathbf{1}_d\| \sqrt{\frac{3s \log(d/\rho)}{t}}.$$

And, integrating over  $\rho$ , we obtain,

$$\mathbb{E}[\|\mathbf{w}_{t,\pi} - \boldsymbol{\omega}\|] \leq O(\|\mathbf{1}_d\| \sqrt{\frac{s \log(d)}{t}}).$$



**High Probability bounds.** For high probability bounds, firstly from Equation (3) and (4),

$$\begin{aligned} \sum_t \mathbb{E}[g_t(\boldsymbol{\theta}_t) | \mathcal{F}_{t-1}] &\leq \sum_t \|\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\| \\ &= \sum_t \|\mathbf{w}_{t,\pi} - \boldsymbol{\omega}\|, \end{aligned}$$

for uniform at random orderings  $\pi$ .

Then, as in above, using Corollary 3 we obtain that for every  $t$ , with probability  $1 - \frac{\rho}{T}$

$$\|\mathbf{w}_{t,\pi} - \boldsymbol{\omega}\| \leq \|\mathbf{1}_d\| \sqrt{\frac{3s \log(dT/\rho)}{t}}.$$

Taking union bound over  $t = 1, \dots, T$ , and summing over  $t$  we obtain that with probability  $1 - \rho$ ,

$$\begin{aligned} \sum_t \mathbb{E}[g_t(\boldsymbol{\theta}_t) | \mathcal{F}_{t-1}] &\leq \sum_t \|\mathbf{w}_{t,\pi} - \boldsymbol{\omega}\| \\ &= O(\|\mathbf{1}_d\| \sqrt{T \log(dT/\rho)}). \end{aligned}$$

Now, using Lemma 16 for dependent random variables  $X_t = g_t(\boldsymbol{\theta}_t)$ , with  $|X_t| = |\boldsymbol{\theta}_t \cdot \mathbf{v}_t^\dagger - h_S(\boldsymbol{\theta}_t)| \leq \|\mathbf{1}_d\|$ , we have,

$$\sum_t g_t(\boldsymbol{\theta}_t) - \sum_t \mathbb{E}[g_t(\boldsymbol{\theta}_t) | \mathcal{F}_{t-1}] \leq O(\|\mathbf{1}_d\| \sqrt{T \log(1/\rho)})$$

with probability at least  $1 - \rho$ .

Combining the above observations, we obtain that with probability  $1 - \rho$ ,

$$\sum_t g_t(\boldsymbol{\theta}_t) \leq O(\|\mathbf{1}_d\| \sqrt{T \log(dT/\rho)}).$$

## Appendix D: Proof of Lemma 7

The offline optimal solution needs to pick  $\mathbf{v}_t^* \in \text{Conv}(X_t)$  to serve request type  $X_t$ , where  $\text{Conv}(X_t)$  denotes the convex hull of set  $X_t$ . Therefore,  $\text{OPT}^\delta$  is defined as

$$\begin{aligned} \text{OPT}^\delta &:= \max_{\{\mathbf{v}_t \in \text{Conv}(X_t)\}} \frac{f(\frac{1}{T} \sum_t \mathbf{v}_t)}{d(\frac{1}{T} \sum_t \mathbf{v}_t, S) \leq \delta} \\ &= \min_{\lambda \geq 0} \max_{\{\mathbf{x} = \frac{1}{T} \sum_t \mathbf{v}_t, \mathbf{v}_t \in \text{Conv}(X_t)\}} \{f(\mathbf{x}) - \lambda d(\mathbf{x}, S) + \delta \lambda\} \\ &= \min_{\lambda \geq 0} \max_{\{\mathbf{x} = \frac{1}{T} \sum_t \mathbf{v}_t, \mathbf{v}_t \in \text{Conv}(X_t)\}} \min_{\|\boldsymbol{\phi}\|_* \leq L, \|\boldsymbol{\theta}\|_* \leq 1} \{f^*(\boldsymbol{\phi}) - \boldsymbol{\phi} \cdot \mathbf{x} - \lambda \boldsymbol{\theta} \cdot \mathbf{x} + \lambda h_S(\boldsymbol{\theta}) + \delta \lambda\} \\ &= \min_{\lambda \geq 0, \|\boldsymbol{\phi}\|_* \leq L, \|\boldsymbol{\theta}\|_* \leq 1} \max_{\{\mathbf{x} = \frac{1}{T} \sum_t \mathbf{v}_t, \mathbf{v}_t \in \text{Conv}(X_t)\}} \{f^*(\boldsymbol{\phi}) - \boldsymbol{\phi} \cdot \mathbf{x} - \lambda \boldsymbol{\theta} \cdot \mathbf{x} + \lambda h_S(\boldsymbol{\theta}) + \delta \lambda\} \\ &= \min_{\lambda \geq 0, \|\boldsymbol{\phi}\|_* \leq L, \|\boldsymbol{\theta}\|_* \leq 1} \left\{ f^*(\boldsymbol{\phi}) + \lambda h_S(\boldsymbol{\theta}) + \frac{1}{T} \sum_{t=1}^T h_{\text{Conv}(X_t)}(-\boldsymbol{\phi} - \lambda \boldsymbol{\theta}) + \delta \lambda \right\} \end{aligned} \tag{14}$$

where, recall that for any convex set  $X$ ,  $h_X(\boldsymbol{\theta})$  was defined as  $h_X(\boldsymbol{\theta}) := \max_{\mathbf{v} \in X} \boldsymbol{\theta} \cdot \mathbf{v}$ . Because a linear function is maximized at a vertex of a convex set,  $h_{\text{Conv}(X_t)}(-\boldsymbol{\phi} - \lambda \boldsymbol{\theta})$  is same as  $h_{X_t}(-\boldsymbol{\phi} - \lambda \boldsymbol{\theta})$ . This allows us to rewrite the expression for  $\text{OPT}^\delta$  as

$$\text{OPT}^\delta = \min_{\lambda \geq 0, \|\boldsymbol{\phi}\|_* \leq L, \|\boldsymbol{\theta}\|_* \leq 1} \left\{ f^*(\boldsymbol{\phi}) + \lambda h_S(\boldsymbol{\theta}) + \frac{1}{T} \sum_{t=1}^T h_{X_t}(-\boldsymbol{\phi} - \lambda \boldsymbol{\theta}) + \delta \lambda \right\} \quad (15)$$

From above, it is clear that  $\text{OPT}^\delta$  is a non-decreasing concave function of  $\delta$ , with gradient as  $\lambda^*(\delta) \geq 0$ , where  $\lambda^*(\delta)$  is the optimal dual variable corresponding to the distance constraint. And,

$$\lim_{\delta \rightarrow 0} \frac{\text{OPT}^\delta - \text{OPT}}{\delta} = \lambda^*$$

where  $\lambda^*$  is the optimal dual variable for  $\text{OPT}$  (i.e., the case of  $\delta = 0$ ). This proves the lemma.

## Appendix E: Proof of Theorem 1

We provide proof of a more detailed theorem statement.

**THEOREM 7.** *Given  $Z$  that satisfies Assumption 1, Algorithm 2 achieves the following regret bounds for online stochastic CP, in RP model:*

$$\begin{aligned} \mathbb{E}[\text{avg-regret}_1(T)] &\leq \frac{(Z+L)}{T} \cdot O(\mathcal{R}(T) + \mathcal{Q}(T)) + O\left(\frac{\mathcal{R}'(T)}{T}\right), \\ \mathbb{E}[\text{avg-regret}_2(T)] &\leq \frac{1}{T} \cdot O(\mathcal{R}(T) + \mathcal{Q}(T)) + \frac{1}{(Z+L)} O\left(\frac{\mathcal{R}'(T)}{T}\right), \end{aligned}$$

where  $\mathcal{Q}(T) = O(\|\mathbf{1}_d\| \sqrt{sT \log(d)})$ ,  $\mathcal{R}'(T)$  is the regret bound for OCO on  $\psi_t(\cdot)$ ,  $\mathcal{R}(T)$  is the regret bound for OCO on  $g_t(\cdot)$ . And,  $s \leq 1$  is the coordinate-wise largest value a vector in  $S$  can take.

Then, substituting OCO regret bounds from Corollary 1 gives the statement of Theorem 1.

*Proof.* Denote by  $(\mathbf{v}_t^*)$  the choice made by the offline optimal solution to satisfy request  $A_t$ . Then,

$$f(\mathbb{E}[\mathbf{v}_t^*]) \geq \text{OPT}, \text{ and } \mathbb{E}[\mathbf{v}_t^*] \in S,$$

where expectation is over  $A_t$  drawn uniformly at random from  $X_1, \dots, X_T$ .

Lemma 18 provides

$$f(\mathbb{E}[\mathbf{v}_t^*]) + \frac{1}{T} \sum_t \mathbb{E}[\psi_t(\boldsymbol{\phi}_t) + 2(Z+L)g_t(\boldsymbol{\theta}_t)] \leq (Z+L) \frac{\mathcal{Q}(T)}{T}$$

where  $\mathcal{Q}(T) = O(\|\mathbf{1}_d\| \sqrt{s \log(d)T})$ . Using Fenchel duality and OCO guarantees, it follows that

$$\begin{aligned} \min_{\|\boldsymbol{\theta}\|_* \leq 1} \frac{1}{T} \sum_t g_t(\boldsymbol{\theta}) &= d\left(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S\right) \\ &\leq \frac{1}{T} \sum_t g_t(\boldsymbol{\theta}_t) + \frac{1}{T} \mathcal{R}(T), \end{aligned}$$

$$\min_{\|\boldsymbol{\phi}\|_* \leq L} \psi_t(\boldsymbol{\phi}) = -f\left(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger\right) \leq \frac{1}{T} \sum_t \psi_t(\boldsymbol{\theta}_t) + \frac{1}{T} \mathcal{R}'(T).$$

Then, using above observations, along with  $f(\mathbb{E}[\mathbf{v}_t^*]) \geq \text{OPT}$ , we obtain

$$\begin{aligned} & \text{OPT} - \mathbb{E}[f(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger)] + 2(Z+L)\mathbb{E}[d(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S)] \\ & \leq \frac{2(Z+L)}{T}(\mathcal{Q}(T) + \mathcal{R}(T)) - \frac{1}{T}\mathcal{R}'(T). \end{aligned}$$

This gives

$$\mathbb{E}[f(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger)] \geq \text{OPT} + 2(Z+L)\mathbb{E}[d(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S)] - \frac{2(Z+L)}{T}(\mathcal{Q}(T) + \mathcal{R}(T)) - \frac{1}{T}\mathcal{R}'(T) \quad (16)$$

Now, we use Assumption 1, to upper bound the reward obtained by the algorithm in terms of OPT and distance from set  $S$ . In particular, we obtain that for  $\delta := \mathbb{E}[d(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S)]$ ,

$$\begin{aligned} \mathbb{E}[f(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger)] & \leq f(\mathbb{E}[\frac{1}{T} \sum_t \mathbf{v}_t^\dagger]) \\ & \leq \text{OPT}^\delta \\ & \leq \text{OPT} + Z\delta \\ & = \text{OPT} + Z \cdot \mathbb{E}[d(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S)]. \end{aligned} \quad (17)$$

Combining the above two inequalities, we obtain

$$\mathbb{E}[d(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S)] \leq \frac{2}{T}(\mathcal{R}(T) + \mathcal{Q}(T)) + \frac{1}{(Z+L)}\mathcal{R}'(T).$$

And, from (16) (using  $\mathbb{E}[d(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger, S)] \geq 0$ ),

$$\mathbb{E}[f(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger)] \geq \text{OPT} - \frac{2(Z+L)}{T} \cdot (\mathcal{R}(T) + \mathcal{Q}(T)) - \frac{\mathcal{R}'(T)}{T}. \quad (18)$$

This gives the theorem statement. □

LEMMA 18.

$$f(\mathbb{E}[\mathbf{v}_t^*]) + \frac{1}{T} \sum_t \mathbb{E}[\psi_t(\phi_t) + 2(Z+L)g_t(\theta_t)] \leq \frac{1}{T}(Z+L)O(\|\mathbf{1}_d\|\sqrt{sT \log(d)}).$$

*Proof.* On the same lines as the proof for BwC, we can prove:

$$\begin{aligned} \psi_t(\phi_t) + 2(Z+L)g_t(\theta_t) & = \phi \cdot \mathbf{v}_t^\dagger - (-f)^*(\phi) + 2(Z+L)(\theta_t \cdot \mathbf{v}_t^\dagger - h_S(\theta_t)) \\ & \leq \phi \cdot \mathbf{v}_t^* - (-f)^*(\phi) + 2(Z+L)(\theta_t \cdot \mathbf{v}_t^* - h_S(\theta_t)). \end{aligned}$$

$$\begin{aligned} \mathbb{E}[\psi_t(\phi_t) + 2(Z+L)g_t(\theta_t) | \mathcal{F}_{t-1}] & \leq \phi_t \cdot \mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - (-f)^*(\phi_t) + 2(Z+L)(\theta_t \cdot \mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - h_S(\theta_t)) \\ & \leq -f(\mathbb{E}[\mathbf{v}_t^*]) + \phi_t \cdot (\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]) + 2(Z+L)\theta_t \cdot (\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]), \end{aligned}$$

where the last inequality uses  $\phi_t \cdot \mathbb{E}[\mathbf{v}_t^*] - (-f)^*(\phi_t) \leq -f(\mathbb{E}[\mathbf{v}_t^*])$  (using Fenchel duality) and  $\theta_t \cdot \mathbb{E}[\mathbf{v}_t^*] - h_S(\theta_t) \leq d(\mathbb{E}[\mathbf{v}_t^*], S) = 0$ . Then, as in proof of Lemma 6,  $\mathbb{E}[\sum_t \|\mathbb{E}[\mathbf{v}_t^* | \mathcal{F}_{t-1}] - \mathbb{E}[\mathbf{v}_t^*]\|]$  can be upper bounded by  $O(\sqrt{\|\mathbf{1}_d\|sT \log(d)})$ . Using this along with observation that  $\|\phi_t\|_* \leq L, \|\theta_t\|_* \leq 1$ , we get the desired lemma statement. □

## Appendix F: Estimating the parameter $Z$

Let  $Z^*$  denote the minimum value of  $Z$  that satisfies the property in Equation (6). As discussed in the proof of Lemma 7,  $Z^* = \lambda^*$ , the value of optimal dual variable corresponding to feasibility constraint. To obtain low regret bounds, ideally we would like to use  $Z = Z^*$  in Algorithm 2, which would provide the minimum possible regret bound of  $O((Z^* + L)\sqrt{\frac{C}{T}})$  in objective according to Theorem 1. The regret in constraints does not depend on  $Z$ . However, in the absence of knowledge of  $Z^*$ , we need to obtain a good enough approximation. Following lemma provides a relaxed condition to be satisfied by  $Z$  in order to obtain the same order of regret bounds, as those obtained with  $Z = Z^*$ .

LEMMA 19. *Assume that  $Z \geq 0$  satisfies the following property, for all  $\delta \geq 3\gamma$  where  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(dT)}{T}}$ ,*

$$\frac{OPT^\delta - OPT^{2\gamma}}{\delta} \leq Z = O(Z^* + L).$$

*Then, Algorithm 2 using such a  $Z$  will achieve an expected regret bound of  $O((Z^* + L)\gamma)$  in objective, and  $O(\gamma)$  in constraints.*

*To compare with Theorem 1, note that  $\gamma = O(\sqrt{\frac{C \log(T)}{T}})$ , therefore, using such a  $Z$  degrades the regret bounds by only an  $O(\sqrt{\log(T)})$  factor.*

*Proof.* Recall that in the proof of Theorem 1, the condition  $OPT^\delta \leq OPT + Z\delta$  was used in the following way. We had the inequality,

$$\begin{aligned} OPT^{\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]} &\geq \mathbb{E}[f(\mathbf{v}_{\text{avg}}^\dagger)] \\ &\geq OPT + 2(Z + L)\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] - \ell(T), \end{aligned} \quad (19)$$

where  $\ell(T) = O((Z + L)\sqrt{\frac{C}{T}})$ . Then, we applied  $OPT^{\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]} \leq OPT + Z\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]$ , to obtain  $OPT + Z\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \geq OPT + 2(Z + L)\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] - \ell(T)$ , yielding  $\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \leq \frac{1}{(Z+L)}O(\ell(T)) = O(\sqrt{\frac{C}{T}})$ .

Now, we will show that it suffices to have  $Z \geq \frac{OPT^\delta - OPT^{2\gamma}}{\delta}$ , for  $\delta > 3\gamma$  to obtain the given regret bounds.

We first bound  $\mathbb{E}[\text{avg-regret}_2(T)] = \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]$ . Starting with Equation 19, observe that if  $\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \leq 3\gamma$ , then the distance is bounded by  $O(\gamma)$  as required anyway, therefore, assume that  $\delta := \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \geq 3\gamma$ . Then, from the given property of  $Z$  we have  $OPT^{\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]} = OPT^\delta \leq OPT^{2\gamma} + Z\delta = OPT^{2\gamma} + Z\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)]$ . Substituting back in Equation (19), we get

$$OPT^{2\gamma} + Z\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \geq OPT + 2(Z + L)\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] - \ell(T)$$

which gives

$$\begin{aligned} (Z + L)\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] &\leq \ell(T) + OPT^{2\gamma} - OPT \\ &\leq \ell(T) + 2Z^*\gamma \\ &= O((Z + L)\gamma) + 2Z^*\gamma \end{aligned}$$

Then, using  $Z = O(Z^* + L)$ , we get

$$\begin{aligned}\mathbb{E}[\text{avg-regret}_2(T)] &= \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] \\ &= O(\gamma) = O\left(\sqrt{\frac{C \log(T)}{T}}\right).\end{aligned}$$

The bound on  $\mathbb{E}[\text{avg-regret}_1(T)]$  depends only on the upper bound on  $Z$  used, and  $Z = O(Z^* + L)$  makes this regret bound to be  $O((Z^* + L)\sqrt{\frac{C}{T}})$ .  $\square$

Next, we provide method for estimating a  $Z$  that satisfies the property stated in Lemma 19. Define

$$\hat{\text{OPT}}^\delta(n) = \max_{\{\mathbf{v}_t \in \text{Conv}(A_t)\}} d\left(\frac{1}{n} \sum_{t=1}^n \mathbf{v}_t, S\right) \leq \delta \quad (20)$$

with  $\hat{\text{OPT}}(n)$  denoting  $\hat{\text{OPT}}^\delta(n)$  for  $\delta = 0$ . We will divide the timeline into phases of size  $1, 1, 2^1, 2^2, \dots, 2^r, \dots$ . Note that phase  $r \geq 2$  consists of  $T_r = 2^{r-2}$  time steps, and there are  $T_r$  time steps before phase  $r$ . The first phase of a single step, we make an arbitrary choice. Then, in every phase  $r \geq 2$ , we will rerun the algorithm, using  $Z$  constructed using observations from the previous  $T_r$  time steps as

$$Z := \frac{(\hat{\text{OPT}}^{4\gamma}(T_r) - \hat{\text{OPT}}^\gamma(T_r))}{\gamma} + 2L \quad (21)$$

with  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(dT_r)}{T_r}}$ .

We prove the following lemma regarding the estimate  $Z$  used in above. Here we use the observation that in RP model, the first  $n$  time steps provide a random sample of observations from the  $T$  observations.

LEMMA 20. *For all  $\rho > 0$  and for all natural numbers  $n$ , let  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(d/\rho)}{n}}$ , and*

$$Z := \frac{(\hat{\text{OPT}}^{4\gamma}(n) - \hat{\text{OPT}}^\gamma(n))}{\gamma} + 2L.$$

*Then, for all  $\delta > 3\gamma$ , with probability  $1 - O(\rho)$ ,*

$$\frac{(\text{OPT}^\delta - \text{OPT}^{2\gamma})}{\delta} \leq Z \leq O(L + Z^*).$$

The proof of above lemma is provided later. We now state the regret bounds for Algorithm 8.

THEOREM 8. *Algorithm 8 has an expected regret of  $\tilde{O}(\sqrt{\frac{C}{T}})$  in the objective and  $(Z^* + L)\tilde{O}(\sqrt{\frac{C}{T}})$  in the constraints.*

*Proof.* For phase  $r \geq 2$ , using  $n = 2^{r-2} = T_r$ , the number of time steps in phase  $r$ , and  $\rho = \frac{1}{T_r^2}$ , from Lemma 20 we obtain that with probability  $1 - O(\frac{1}{T_r^2})$ ,  $Z$  available to phase  $r$  satisfies the property required by Lemma 19 (with  $T$  substituted by  $T_r$ ), which gives the following regret bounds for phase  $r$ : let  $\mathbf{v}_{\text{avg}}^\dagger(r)$  be the average of played vectors in the  $T_r$  time steps of phase  $r$ . Let  $\mathcal{F}_{r-1}$  denote the

history till phase  $r - 1$ . Then, with probability  $1 - O(\frac{1}{T^2})$  the history  $\mathcal{F}_{r-1}$  is such that in phase  $r$  the regret in distance is bounded by  $\mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger(r), S) | \mathcal{F}_{r-1}] \leq \tilde{O}(\sqrt{\frac{C}{T_r}})$ . With remaining probability  $O(\frac{1}{T^2})$ , the distance can be at most  $T_r \|\mathbf{1}_d\|$ . Let  $\mathbf{v}_{\text{avg}}^\dagger$  denote the average of played vectors from the entire period of  $T$  time steps. Then, we get that total regret,

$$\begin{aligned} \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger, S)] &\leq \frac{\|\mathbf{1}_d\|}{T} + \sum_{r=2}^{\log(T)+1} \frac{T_r}{T} \mathbb{E}[d(\mathbf{v}_{\text{avg}}^\dagger(r), S)] \\ &\leq \frac{\|\mathbf{1}_d\|}{T} + \sum_{r=2}^{\log(T)+1} \frac{T_r}{T} \tilde{O}(\sqrt{\frac{C}{T_r}} + \frac{T_r \|\mathbf{1}_d\|}{T_r^2}) \\ &= \tilde{O}(\sqrt{\frac{C}{T}}). \end{aligned}$$

Similarly, we obtain bounds on regret in the objective,

$$\begin{aligned} \text{OPT} - \mathbb{E}[f(\mathbf{v}_{\text{avg}}^\dagger)] &\leq \frac{1}{T} + \sum_{r=2}^{\log(T)+1} \frac{T_r}{T} (\text{OPT} - \mathbb{E}[f(\mathbf{v}_{\text{avg}}^\dagger(r))]) \\ &\leq \frac{1}{T} + \sum_{r=2}^{\log(T)+1} \frac{T_r}{T} (Z^* + L) \tilde{O}(\sqrt{\frac{C}{T_r}} + \frac{T_r \|\mathbf{1}_d\|}{T_r^2}) \\ &= (Z^* + L) \tilde{O}(\sqrt{\frac{C}{T}}). \end{aligned}$$

□

*Proof of Lemma 20.* From Lemma 7,  $\text{OPT}^\delta$  is concave in  $\delta$ , therefore, for all  $\delta > 3\gamma$

$$\begin{aligned} \frac{(\text{OPT}^\delta - \text{OPT}^{2\gamma})}{\delta} &\leq \frac{(\text{OPT}^\delta - \text{OPT}^{2\gamma})}{\delta - 2\gamma} \\ &\leq \frac{(\text{OPT}^{3\gamma} - \text{OPT}^{2\gamma})}{\gamma}. \end{aligned}$$

So, it suffices to prove that

$$\frac{(\text{OPT}^{3\gamma} - \text{OPT}^{2\gamma})}{\gamma} \leq Z \leq O(L + Z^*).$$

In Lemma 22 and Lemma 23, we prove that for every  $\delta \geq \gamma$ , with probability  $1 - O(\rho)$

$$\begin{aligned} \hat{\text{OPT}}^\delta + L\gamma &\geq \text{OPT}^{\delta-\gamma}, \\ \text{OPT}^\delta + L\gamma &\geq \hat{\text{OPT}}^{\delta-\gamma} \end{aligned} \tag{22}$$

Using above for  $\delta = 4\gamma$ , and  $\delta = 2\gamma$ , respectively, we get

$$\begin{aligned} Z &:= \frac{(\hat{\text{OPT}}^{4\gamma}(n) - \hat{\text{OPT}}^{2\gamma}(n))}{\gamma} + 2L \\ &\geq \frac{(\text{OPT}^{3\gamma} - \text{OPT}^{2\gamma})}{\gamma}. \end{aligned}$$

In Lemma 24, we prove that for any  $\delta \geq \gamma$ ,

$$\hat{\text{OPT}}^\delta \leq \text{OPT} + O(\delta(Z^* + L)) \tag{23}$$

Using this along with  $\hat{\text{OPT}}^\gamma \geq \text{OPT} - L\gamma$  from the first inequality in Equation (22), we get

$$\begin{aligned} Z &= \frac{\hat{\text{OPT}}^{4\gamma} - \hat{\text{OPT}}^\gamma}{\gamma} \\ &\leq \frac{(\text{OPT} + 4\gamma O(Z^* + L)) - (\text{OPT} - L\gamma)}{\gamma} \\ &= O(Z^* + L). \end{aligned}$$

This completes the proof.  $\square$

LEMMA 21. *Given fixed  $\{\mathbf{v}_t\}_{t=1}^T$ , and a vector  $\boldsymbol{\mu}$ , for all  $\rho > 0$  and  $n \in [T]$ , let  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(d/\rho)}{n}}$ . Then for a uniformly random permutation over  $1, \dots, T$ , with probability  $1 - O(\rho)$ , the following holds for the first  $n$  time steps.*

$$\begin{aligned} \left\| \frac{1}{n} \sum_{t=1}^n \mathbf{v}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t \right\| &\leq \gamma, \\ \left| \frac{1}{n} \sum_{t=1}^n h_{A_t}(\boldsymbol{\mu}) - \frac{1}{T} \sum_{t=1}^T h_{A_t}(\boldsymbol{\mu}) \right| &\leq \gamma \|\boldsymbol{\mu}\|_*. \end{aligned}$$

*Proof.* The first inequality is obtained by simple application of Chernoff-Hoeffding bounds (Lemma 14) for every coordinate  $v_{t,j}$ , which gives

$$\left| \frac{1}{n} \sum_{t=1}^n v_{t,j} - \frac{1}{T} \sum_{t=1}^T v_{t,j} \right| \leq \sqrt{\frac{\log(d/\rho)}{n}},$$

with probability  $1 - O(\rho/d)$ . Then taking union bound over the  $d$  coordinates, we get the required inequality.

The second inequality follows using Chernoff-Hoeffding bounds (Lemma 14) for bounded random variables  $Y_t = h_{A_t}(\boldsymbol{\mu})$ , where  $|Y_t| = |h_{A_t}(\boldsymbol{\mu})| \leq \|\boldsymbol{\mu}\|_* \cdot \|\mathbf{1}_d\|$  (from the definition of the dual norm). This gives with probability  $1 - O(\rho)$ ,

$$\begin{aligned} \left| \frac{1}{n} \sum_{t=1}^n h_{A_t}(\boldsymbol{\mu}) - \frac{1}{T} \sum_{t=1}^T h_{A_t}(\boldsymbol{\mu}) \right| &= \left| \frac{1}{n} \sum_{t=1}^n (Y_t - \mathbb{E}[Y_t]) \right| \\ &\leq (\|\boldsymbol{\mu}\|_* \cdot \|\mathbf{1}_d\|) \sqrt{\frac{\log(1/\rho)}{n}} \\ &\leq \|\boldsymbol{\mu}\|_* \gamma. \end{aligned}$$

$\square$

LEMMA 22. *For all  $\rho > 0$  and  $n \in [T]$ , let  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(d/\rho)}{n}}$ . For all  $\delta \geq \gamma$ , with probability  $1 - O(\rho)$ ,*

$$\hat{\text{OPT}}^\delta(n) \geq \text{OPT}^{\delta-\gamma} - L\gamma.$$

*Proof.* To prove  $\hat{\text{OPT}}^\delta(n) \geq \text{OPT}^{\delta-\gamma} - L\gamma$ , we prove that there exists a feasible primal solution of  $\hat{\text{OPT}}^\delta(n)$  that is at most  $\gamma$  distance from the optimal primal solution of  $\text{OPT}^{\delta-\gamma}$ . Then, the lemma follows from the  $L$ -Lipschitz property of  $f$ .

Let  $\{\mathbf{v}_t\}_{t=1}^T$  be the optimal primal solution for  $\text{OPT}^{\delta-\gamma}$ , so that  $d(\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t, S) \leq \delta - \gamma$ . Then,

$$\begin{aligned} d\left(\frac{1}{n} \sum_{t=1}^n \mathbf{v}_t, S\right) &\leq \left\| \frac{1}{n} \sum_{t=1}^n \mathbf{v}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t \right\| + d\left(\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t, S\right) \\ &\leq \gamma + (\delta - \gamma) = \delta, \end{aligned}$$

where we used the concentration bounds from Lemma 21 to bound  $\left\| \frac{1}{n} \sum_{t=1}^n \mathbf{v}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t \right\|$  by  $\gamma$ . Therefore,  $\{\mathbf{v}_t\}_{t=1}^n$  is a primal feasible solution of  $\hat{\text{OPT}}^\delta(n)$  with objective value  $f(\frac{1}{n} \sum_{t=1}^n \mathbf{v}_t) \geq f(\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t) - L\left\| \frac{1}{n} \sum_{t=1}^n \mathbf{v}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{v}_t \right\| \geq f(\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t) - L\gamma = \text{OPT}^{\delta-\gamma} - L\gamma$ . Therefore,  $\hat{\text{OPT}}^\delta \geq f(\frac{1}{n} \sum_{t=1}^n \mathbf{v}_t) \geq \text{OPT}^{\delta-\gamma} - L\gamma$ .  $\square$

---

**Algorithm 8** Algorithm for online CP with  $Z$  estimation

---

Choose any option in the first step.

**for all** phases  $r = 2, \dots, \log(T) + 1$  **do**

    COMPUTE  $Z$  using observations in steps 1 to  $T_r = 2^{r-2}$  as

$$Z = \frac{(\hat{\text{OPT}}^{4\gamma}(T_r) - \hat{\text{OPT}}^\gamma(T_r))}{\gamma} + 2L$$

    with  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(dT_r)}{T_r}}$ .

    Run Algorithm 2 for  $T_r$  steps  $t = \{T_r + 1, \dots, 2T_r\}$  of phase  $r$  using  $Z$  as computed above.

**end for**

---

LEMMA 23. For all  $\rho > 0$  and  $n \in [T]$ , let  $\gamma = \|\mathbf{1}_d\| \sqrt{\frac{\log(d/\rho)}{n}}$ . For all  $\delta \geq \gamma$ , with probability  $1 - O(\rho)$ ,

$$\text{OPT}^\delta + L\gamma \geq \hat{\text{OPT}}^{\delta-\gamma}(n).$$

*Proof.* Define  $S^\delta$  as the set  $\{\mathbf{v} : d(\mathbf{v}, S) \leq \delta\}$ . Then, using the derivation in Equation (15), we have that

$$\text{OPT}^\delta = \min_{\lambda \geq 0, \|\phi\|_* \leq L, \|\theta\|_* \leq 1} \left\{ f^*(\phi) + \lambda h_{S^\delta}(\theta) + \frac{1}{T} \sum_{t=1}^T h_{A_t}(-\phi - \lambda\theta) \right\}.$$

Let  $\lambda^*, \theta^*, \phi^*$  be the optimal dual solutions in above. Then,

$$\begin{aligned} \hat{\text{OPT}}^{\delta-\gamma}(n) &= \min_{\lambda \geq 0, \|\phi\|_* \leq L, \|\theta\|_* \leq 1} \left\{ f^*(\phi) + \lambda h_{S^{\delta-\gamma}}(\theta) + \frac{1}{n} \sum_{t=1}^n h_{A_t}(-\phi - \lambda\theta) \right\} \\ &\leq f^*(\phi^*) + \lambda^* h_{S^{\delta-\gamma}}(\theta^*) + \frac{1}{n} \sum_{t=1}^n h_{A_t}(-\phi^* - \lambda^*\theta^*) \end{aligned}$$

Now, using concentration bounds from Lemma 21 for the sum of  $h_{A_t}$ 's, we obtain,

$$\hat{\text{OPT}}^{\delta-\gamma}(n) \leq f^*(\phi^*) + \lambda^* h_{S^{\delta-\gamma}}(\theta^*) + \frac{1}{T} \sum_{t=1}^T h_{A_t}(-\phi^* - \lambda^*\theta^*) + \gamma(\lambda^* \|\theta^*\|_* + \|\phi^*\|_*).$$



Now, observe that for any  $\theta$ ,  $h_{S^\delta}(\theta) \geq h_{S^{\delta-\gamma}}(\theta) + \gamma\|\theta\|_*$ . To see this, let  $\mathbf{v}$  be the maximizer in the definition of  $h_{S^{\delta-\gamma}}$ , i.e.,  $\mathbf{v} = \arg \max_{\mathbf{u} \in S^{\delta-\gamma}} \mathbf{u} \cdot \theta$ . Then consider  $\mathbf{v}' = \mathbf{v} + \gamma \frac{\theta}{\|\theta\|}$ . We have that  $\|\mathbf{v}' - \mathbf{v}\| = \gamma$ , so that  $\mathbf{v} \in S^{\delta-\gamma}$  implies that  $\mathbf{v}' \in S^\delta$ . Therefore  $h_{S^\delta}(\theta) \geq \mathbf{v}' \cdot \theta = \mathbf{v} \cdot \theta + \gamma\|\theta\|_* = h_{S^{\delta-\gamma}}(\theta) + \gamma\|\theta\|_*$ . Substituting, we get,

$$\begin{aligned} \widehat{\text{OPT}}^{\delta-\gamma}(n) &\leq f^*(\phi^*) + \lambda^* h_{S^\delta}(\theta^*) - \gamma\lambda^*\|\theta^*\|_* + \frac{1}{T} \sum_{t=1}^T h_{A_t}(-\phi^* - \lambda^*\theta^*) + \gamma(\lambda^*\|\theta^*\|_* + \|\phi^*\|_*) \\ &= \text{OPT}^\delta + \gamma\|\phi^*\|_* \\ &\leq \text{OPT}^\delta + \gamma L \end{aligned}$$

□

LEMMA 24. For all  $\delta \geq \gamma$ , with probability  $1 - O(\rho)$ ,

$$O\hat{P}T^\delta(n) \leq OPT + 2\delta(L + Z^*)$$

*Proof.* Using the derivations in Equation (15),

$$\widehat{\text{OPT}}^\delta(n) = \min_{\lambda \geq 0, \|\phi\|_* \leq L, \|\theta\|_* \leq 1} \left\{ f^*(\phi) + \lambda h_S(\theta) + \frac{1}{n} \sum_{t=1}^n h_{A_t}(-\phi - \lambda\theta) + \delta\lambda \right\},$$

Let  $\lambda^*, \phi^*, \theta^*$  denote the optimal dual solution for OPT, then,

$$\widehat{\text{OPT}}^\delta(n) \leq f^*(\phi^*) + \lambda h_S(\theta^*) + \frac{1}{n} \sum_{t=1}^n h_{A_t}(-\phi^* - \lambda^*\theta^*) + \delta\lambda^*$$

Now, using concentration bounds from Lemma 21 for the sum of  $h_{A_t}$ 's, we obtain,

$$\begin{aligned} O\hat{P}T^\delta(n) &\leq f^*(\phi^*) + \lambda h_S(\theta^*) + \frac{1}{T} \sum_{t=1}^T h_{A_t}(-\phi^* - \lambda^*\theta^*) + \gamma(\lambda^*\|\theta^*\|_* + \|\phi^*\|_*) + \delta\lambda^* \\ &= \text{OPT} + \gamma(\lambda^*\|\theta^*\|_* + \|\phi^*\|_*) + \delta\lambda^* \\ &\leq \text{OPT} + (L + \lambda^*)\gamma + \delta\lambda^* \\ &\leq \text{OPT} + 2(L + \lambda^*)\delta \\ &= \text{OPT} + 2\delta(L + Z^*) \end{aligned}$$

□

## Appendix G: Proof of Lemma 11

Given an instance of the online packing problem, recall that  $(r_t^*, \mathbf{v}_t^*)$  denotes the optimal offline solution. Then  $\text{OPT}_{\text{SUM}} = \sum_{t=1}^T r_t^*$ , and  $\sum_{t=1}^T \mathbf{v}_t^* \leq B\mathbf{1}$ . Given  $\rho > 0$ , let  $\eta = \sqrt{3 \log(\frac{d+2}{\rho})}$ . Let the given random subset of  $\delta$  fraction of requests be  $\Gamma$ . Define  $O\hat{P}T$  to be  $1/\delta$  times the optimum value of the following scaled optimization problem: pick  $(r_t^\dagger, \mathbf{v}_t^\dagger)$  for each  $t \in \Gamma$ , to maximize the total reward  $\sum_{t \in \Gamma} r_t^\dagger$  such that  $\sum_{t \in \Gamma} \mathbf{v}_t^\dagger \leq (\delta B + \eta\sqrt{\delta B})\mathbf{1}$ .

The bounds we need on  $\hat{\text{OPT}}$  follow from considering the optimal primal and dual solutions to the given packing problem restricted to the sample and using Corollary 3 to bound their values on the sample. Applying Corollary 3 to the set of  $r_t^*$  for all  $t \in [T]$  we get that with probability at least  $1 - \rho/(d+2)$ ,

$$\begin{aligned} \sum_{t \in \Gamma} r_t^* &\geq \delta \text{OPT}_{\text{SUM}} - \sqrt{3\delta \text{OPT}_{\text{SUM}} \log\left(\frac{d+2}{\rho}\right)} \\ &= \delta \text{OPT}_{\text{SUM}} - \eta \sqrt{\delta \text{OPT}_{\text{SUM}}}. \end{aligned}$$

Similarly, applying Corollary 3 to each co-ordinate of the set of  $\mathbf{v}_t^*$ s, and taking a union bound, we get that with probability at least  $1 - \rho d/(d+2)$ ,

$$\begin{aligned} \sum_{t \in \Gamma} \mathbf{v}_t^* &\leq (\delta B + \sqrt{3\delta B \log\left(\frac{d+2}{\rho}\right)}) \mathbf{1} \\ &= (\delta B + \eta \sqrt{\delta B}) \mathbf{1}. \end{aligned}$$

Therefore with probability  $1 - \rho(d+1)/(d+2)$  both the inequalities above hold and  $(r_t^*, \mathbf{v}_t^*)_{t \in \Gamma}$  is a feasible solution to the scaled optimization problem used to define  $\hat{\text{OPT}}$ . Hence

$$\delta \hat{\text{OPT}} \geq \sum_{t \in \Gamma} r_t^* \geq \delta \text{OPT}_{\text{SUM}} - \eta \sqrt{\delta \text{OPT}_{\text{SUM}}}$$

and the first bound on  $\hat{\text{OPT}}$  follows from dividing the above inequality throughout by  $\delta$ . For the second bound, we need to consider the dual of the packing problem. The packing problem has the following natural LP relaxation. (The dual LP follows.)

$$\begin{aligned} \max \quad & \sum_{t=1}^T \sum_{\mathbf{v} \in A_t} r(\mathbf{v}) x_{t,\mathbf{v}} \\ \text{s.t.} \quad & \forall t, \sum_{\mathbf{v} \in A_t} x_{t,\mathbf{v}} \leq 1 \\ & \sum_{t=1}^T \sum_{\mathbf{v} \in A_t} \mathbf{v} x_{t,\mathbf{v}} \leq B \mathbf{1}. \end{aligned}$$

$$\begin{aligned} \min \quad & \sum_{t=1}^T \beta_t + B \boldsymbol{\theta} \cdot \mathbf{1} \\ \text{s.t.} \quad & \forall t, \forall \mathbf{v} \in A_t, \beta_t \geq r(\mathbf{v}) - \mathbf{v} \cdot \boldsymbol{\theta}, \\ & \forall t, \beta_t \geq 0, \boldsymbol{\theta} \geq 0. \end{aligned}$$

First of all, we ignore the integrality gap and assume that the value of the optimal dual (and primal) solution is equal to the optimal value  $\text{OPT}_{\text{SUM}}$  for the offline packing problem. Let  $(\beta_t^*)_{t=1}^T, (\theta_j^*)_{j=1}^d$  be the optimal dual solution for the given instance, and  $\text{OPT}_{\text{SUM}} = \sum_t \beta_t^* + \sum_j B \theta_j^*$ . It can be shown that  $\beta_t^* \in [0, 1]$  for all  $t$ : all the constraints involving  $\beta_t$  are of the form  $\beta_t \geq (\cdot)$  so at least one of these

constraints is tight for the optimal solution. Also for each of these constraints, the RHS is at most 1, and one of the constraints is  $\beta_t \geq 0$ . Further note that these constraints are local, i.e., they only depend on the request indexed by  $t$ . This means that  $(\beta_t^*)_{t \in \Gamma}, (\theta_j^*)_{j=1}^d$  is a feasible solution to the dual of the scaled optimization problem. The objective value of this solution to this dual is

$$\sum_{t \in \Gamma} \beta_t^* + \sum_j (\delta B + \eta \sqrt{\delta B}) \theta_j^* \geq \delta \text{OPT}.$$

Using Corollary 3 on the set of  $\beta_t^*$ s, we get that with probability at least  $1 - \rho/(d+2)$ ,

$$\begin{aligned} \sum_{t \in \Gamma} \beta_t^* &\leq \delta \sum_{t=1}^T \beta_t^* + \sqrt{3\delta \text{OPT}_{\text{SUM}} \log\left(\frac{d+2}{\rho}\right)} \\ &= \delta \sum_{t=1}^T \beta_t^* + \eta \sqrt{\delta \text{OPT}_{\text{SUM}}}. \end{aligned}$$

Putting the two inequalities above together,

$$\begin{aligned} \frac{\delta \text{OPT}}{1 + \eta/\sqrt{\delta B}} &\leq \sum_{t \in \Gamma} \beta_t^* + \delta \sum_j B \theta_j^* \\ &\leq \delta \left( \sum_{t=1}^T \beta_t^* + \sum_j B \theta_j^* \right) + \eta \sqrt{\delta \text{OPT}_{\text{SUM}}} \\ &= \delta \text{OPT}_{\text{SUM}} + \eta \sqrt{\delta \text{OPT}_{\text{SUM}}}. \end{aligned}$$

The lemma follows by taking the union bound over the probabilities for the two inequalities as required. Finally, we ignored the integrality gap, but it is easy to show that this gap is at most  $1 - \frac{1}{B}$ , which can be absorbed in the  $1 + \eta/\sqrt{\delta B}$  factor.

## Endnotes

1. In online learning, the objective value is the sum of reward in every step, which scales with  $T$ , and the regret typically scales with  $\sqrt{T}$ . But in our formulation, the objective  $f(\frac{1}{T} \sum_t \mathbf{v}_t^\dagger)$  is defined over average observations, therefore, to be consistent with the popular terminology, we call our regret ‘average regret’.

2. Note that such a stopping rule does not make sense for a general  $S$ . If  $S$  is downwards closed, then one can consider similar stopping rules in those cases as well.

## References

Abernethy, Jacob, Peter L. Bartlett, Elad Hazan. 2011. Blackwell approachability and low-regret learning are equivalent. *COLT*.

- Aggarwal, Gagan, Gagan Goel, Chinmay Karande, Aranyak Mehta. 2011. Online vertex-weighted bipartite matching and single-bid budgeted allocations. *SODA*.
- Agrawal, S., Z. Wang, Y. Ye. 2014. A dynamic near-optimal algorithm for online linear programming. *Operations Research* **62** 876 – 890.
- Agrawal, Shipra, Nikhil R. Devanur. 2014. Bandits with concave rewards and convex knapsacks. *Proceedings of the Fifteenth ACM Conference on Economics and Computation*. EC '14.
- Arora, Sanjeev, Elad Hazan, Satyen Kale. 2012. The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing* **8**(6) 121–164.
- Babaioff, Moshe, Shaddin Dughmi, Robert Kleinberg, Aleksandrs Slivkins. 2012. Dynamic pricing with limited supply. *EC*.
- Badanidiyuru, Ashwinkumar, Robert Kleinberg, Aleksandrs Slivkins. 2013. Bandits with knapsacks. *FOCS*. 207–216.
- Bahmani, Bahman, Michael Kapralov. 2010. Improved bounds for online stochastic matching. *ESA*. 170–181.
- Blackwell, David. 1956. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics* **6**(1) 1–8.
- Buchbinder, Niv, Kamal Jain, Joseph Seffi Naor. 2007. Online primal-dual algorithms for maximizing ad-auctions revenue. *Proceedings of the 15th Annual European Conference on Algorithms*. ESA'07.
- Chakrabarti, Deepayan, Erik Vee. 2012. Traffic shaping to optimize ad delivery. *Proceedings of the 13th ACM Conference on Electronic Commerce*. EC '12.
- Chen, Peiji, Wenjing Ma, Srinath Mandalapu, Chandrashekhhar Nagarjan, Jayavel Shanmugasundaram, Sergei Vassilvitskii, Erik Vee, Manfai Yu, Jason Zien. 2012. Ad serving using a compact allocation plan. *Proceedings of the 13th ACM Conference on Electronic Commerce*. EC '12.
- Chen, Xiao, Zizhuo Wang. 2013. A near-optimal dynamic learning algorithm for online matching problems with concave returns. <http://arxiv.org/abs/1307.5934>.
- Chen, Ye, Pavel Berkhin, Bo Anderson, Nikhil R. Devanur. 2011. Real-time bidding algorithms for performance-based display ad allocation. *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '11.

- Devanur, Nikhil R., Thomas P. Hayes. 2009. The adwords problem: online keyword matching with budgeted bidders under random permutations. *EC*.
- Devanur, Nikhil R., Zhiyi Huang, Nitish Korula, Vahab S. Mirrokni, Qiqi Yan. 2013. Whole-page optimization and submodular welfare maximization with online bidders. *Proceedings of the Fourteenth ACM Conference on Electronic Commerce. EC '13*.
- Devanur, Nikhil R., Kamal Jain. 2012. Online matching with concave returns. *Proceedings of the Forty-fourth Annual ACM Symposium on Theory of Computing. STOC '12*.
- Devanur, Nikhil R., Kamal Jain, Balasubramanian Sivan, Christopher A. Wilkens. 2011a. Near optimal online algorithms and fast approximation algorithms for resource allocation problems. *EC*.
- Devanur, Nikhil R., Kamal Jain, Balasubramanian Sivan, Christopher A. Wilkens. 2011b. Near optimal online algorithms and fast approximation algorithms for resource allocation problems. Full version, accessible from <http://research.microsoft.com/en-us/um/people/bsivan/>.
- Feldman, J., N. Korula, V. Mirrokni, S. Muthukrishnan, M. Pal. 2009a. Online ad assignment with free disposal. *WINE*.
- Feldman, Jon, Monika Henzinger, Nitish Korula, Vahab S. Mirrokni, Cliff Stein. 2010a. Online stochastic packing applied to display ad allocation. *Proceedings of the 18th Annual European Conference on Algorithms: Part I. ESA'10*.
- Feldman, Jon, Monika Henzinger, Nitish Korula, Vahab S. Mirrokni, Clifford Stein. 2010b. Online stochastic ad allocation: Efficiency and fairness. *CoRR* **abs/1001.5076**.
- Feldman, Jon, Aranyak Mehta, Vahab Mirrokni, S. Muthukrishnan. 2009b. Online stochastic matching: Beating  $1-1/e$ . *FOCS '09: Proceedings of the 2009 50th Annual IEEE Symposium on Foundations of Computer Science*.
- Ghosh, Arpita, Randolph Preston McAfee, Kishore Papineni, Sergei Vassilvitskii. 2009. Bidding for representative allocations for display advertising. *WINE*.
- Goel, Gagan, Aranyak Mehta. 2008. Online budgeted matching in random input models with applications to adwords. *SODA '08: Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*.
- Gupta, Anupam, Marco Molinaro. 2014. How the Experts Algorithm Can Help Solve LPs Online. *Algorithms - ESA 2014, Lecture Notes in Computer Science* **8737** 517–529.

- Hazan, Elad, Amit Agarwal, Satyen Kale. 2007. Logarithmic regret algorithms for online convex optimization. *Mach. Learn.* **69**(2-3).
- Hoeffding, Wassily. 1963. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association* **58**(301) 13–30.
- Kakade, Sham M., Shai Shalev-Shwartz, Ambuj Tewari. 2009. On the duality of strong convexity and strong smoothness: Learning applications and matrix regularization. Tech. rep., Toyota Technological Institute - Chicago, USA. <http://ttic.uchicago.edu/~shai/papers/KakadeShalevTewari09.pdf>.
- Karande, Chinmay, Aranyak Mehta, Ramakrishnan Srikant. 2013. Optimizing budget constrained spend in search advertising. *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining. WSDM '13*.
- Karande, Chinmay, Aranyak Mehta, Pushkar Tripathi. 2011. Online bipartite matching with unknown distributions. *STOC*.
- Karp, R. M., U. V. Vazirani, V. V. Vazirani. 1990. An optimal algorithm for on-line bipartite matching. *Proceedings of the Twenty-second Annual ACM Symposium on Theory of Computing. STOC '90*.
- Kesselheim, Thomas, Andreas Tönnis, Klaus Radke, Berthold Vöcking. 2014. Primal beats dual on online packing LPs in the random-order model. *STOC*.
- Kleinberg, R. 2005. A multiple-choice secretary algorithm with applications to online auctions. *Proceedings of the 16th Annual ACM-SIAM Symposium on Discrete algorithms.* 630–631.
- Kleinberg, Robert, Alex Slivkins, Eli Upfal. 2008. Multi-armed bandits in metric spaces. *STOC*.
- Mahdian, Mohammad, Hamid Nazerzadeh, Amin Saberi. 2012. Online optimization with uncertain information. *ACM Trans. Algorithms* **8**(1).
- Mahdian, Mohammad, Qiqi Yan. 2011. Online bipartite matching with random arrivals: an approach based on strongly factor-revealing LPs. *STOC*.
- Manshadi, Vahideh, Shayan Gharan, Amin Saberi. 2011. Online stochastic matching: Online actions based on offline statistics. *SODA*.
- Mehta, Aranyak, Amin Saberi, Umesh V. Vazirani, Vijay V. Vazirani. 2007. Adwords and generalized online matching. *J. ACM* **54**(5).

Mirrokn, Vahab S., Shayan Oveis Gharan, Morteza Zadimoghaddam. 2012. Simultaneous approximations for adversarial and stochastic online budgeted allocation. *Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA '12.

Shalev-Shwartz, Shai. 2012. Online learning and online convex optimization. *Foundations and Trends in Machine Learning* **4**(2) 107–194.

Vee, Erik, Sergei Vassilvitskii, Jayavel Shanmugasundaram. 2010. Optimal online assignment with forecasts. *EC '10: Proceedings of the 11th ACM conference on Electronic commerce*.