This article was downloaded by: [132.174.252.179] On: 06 June 2021, At: 22:28 Publisher: Institute for Operations Research and the Management Sciences (INFORMS)

INFORMS is located in Maryland, USA



## **Operations Research**

Publication details, including instructions for authors and subscription information: http://pubsonline.informs.org

## Online Allocation and Pricing: Constant Regret via Bellman Inequalities

Alberto Vera, Siddhartha Banerjee, Itai Gurvich

### To cite this article:

Alberto Vera, Siddhartha Banerjee, Itai Gurvich (2021) Online Allocation and Pricing: Constant Regret via Bellman Inequalities. **Operations Research** 

Published online in Articles in Advance 26 Mar 2021

. <a href="https://doi.org/10.1287/opre.2020.2061">https://doi.org/10.1287/opre.2020.2061</a>

Full terms and conditions of use: https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-**Conditions** 

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or quarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a quarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2021, INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit http://www.informs.org

Articles in Advance, pp. 1-20 ISSN 0030-364X (print), ISSN 1526-5463 (online)

### **Crosscutting Areas**

## Online Allocation and Pricing: Constant Regret via **Bellman Inequalities**

Alberto Vera, Siddhartha Banerjee, Itai Gurvich

a School of Operations Research and Information Engineering, Cornell University, Ithaca, New York 14853; b Kellogg School of Management, Northwestern University, Evanston, Illinois 60208

Contact: aav39@cornell.edu, Dhttps://orcid.org/0000-0003-0569-8623 (AV); sbanerjee@cornell.edu, 📵 https://orcid.org/0000-0002-8954-4578 (SB); gurvich@cornell.edu, 📵 https://orcid.org/0000-0001-9746-7755 (IG)

Received: June 17, 2019 Revised: May 30, 2020 Accepted: June 11, 2020

Published Online in Articles in Advance:

March 26, 2021

Subject classifications: dynamic programming/ optimal control; applications

Area of Review: Revenue Management and Market Analytics

https://doi.org/10.1287/opre.2020.2061

Copyright: © 2021 INFORMS

**Abstract.** We develop a framework for designing simple and efficient policies for a family of online allocation and pricing problems that includes online packing, budget-constrained probing, dynamic pricing, and online contextual bandits with knapsacks. In each case, we evaluate the performance of our policies in terms of their regret (i.e., additive gap) relative to an offline controller that is endowed with more information than the online controller. Our framework is based on Bellman inequalities, which decompose the loss of an algorithm into two distinct sources of error: (1) arising from computational tractability issues, and (2) arising from estimation/prediction of random trajectories. Balancing these errors guides the choice of benchmarks, and leads to policies that are both tractable and have strong performance guarantees. In particular, in all our examples, we demonstrate constant-regret policies that only require resolving a linear program in each period, followed by a simple greedy action-selection rule; thus, our policies are practical as well as provably near optimal.

Funding: This work was supported by the U.S. Department of Defense [Grants STTR A18B-T007 and W911NF-20-C-0008], the National Science Foundation [Grants CNS-1955997, DMS-1839346, and ECCS-1847393], and the Army Research Laboratory [Grant W911NF-17-1-0094].

Supplemental Material: The online appendices are available at https://doi.org/10.1287/opre.2020.2061.

Keywords: stochastic optimization • approximate dynamic programming • online resource allocation • dynamic pricing • online packing • network revenue management

#### 1. Introduction

Online decision making under uncertainty is widely studied across a variety of fields, including operations research, control, and computer science. A canonical framework for such problems is that of Markov decision processes (MDPs), with associated use of stochastic dynamic programming for designing policies. In complex settings, however, such approaches suffer from the known curse of dimensionality; moreover, they also fail to provide insights into structural properties of the problem: the performance of heuristics, dependence on distributional information, and so on.

These challenges have inspired an alternate approach to designing approximate policies for MDPs based on the use of benchmarks—proxies for the value function that provide bounds for the optimal policy and guide the design of heuristics. The performance of any policy can be quantified by its additive loss, or regret, relative to any such benchmark; this consequently also bounds the additive optimality gap, that is, performance against the optimal policy.

In this work, we develop new policies for online resource-allocation problems: settings where a finite set of resources is dynamically allocated to arriving requests, with associated constraints and rewards/ costs. Our baseline problem is the online stochastic knapsack problem (henceforth OnlineKnapsack): a controller has initial inventory B, and requests arrive sequentially over horizon T. Each request has a random type corresponding to a resource requirementreward pair. Requests are generated from a known stochastic process, and are revealed upon arrival; the controller must then decide whether to accept/reject each request in order to maximize rewards while satisfying budget constraints. We then consider three variants of this basic setting: (1) online probing, (2) dynamic pricing, and (3) contextual bandits with knapsacks. These are widely studied problems, each of which augments the baseline Online Knapsack with additional constraints/controls. The formal models for these settings are presented in Section 2.

Instead of solving each problem in an ad hoc manner, however, our policies are all derived from a single underlying framework. In particular, our results can be summarized as follows:

Meta-Theorem. Given an online allocation problem, we identify an appropriate offline benchmark, and give a simple online policy—based on solving a tractable optimization problem in each period—that gets constant regret compared with the benchmark (and thus, compared with the optimal policy).

In more detail, our approach is based on adaptively constructing a benchmark that has additional (but not necessarily full) information about future randomness. Next, in the spirit of online primal-dual methods, we use our benchmark to construct a feasible online policy. The centerpiece of our approach are the Bellman inequalities, which characterize what benchmarks are feasible and decompose the regret of an online policy into two distinct terms. The first, which we call the Bellman loss, arises from computational considerations, specifically, from requiring that the benchmark is tractable (instead of a dynamic program that may be intractable). The second, which we call the information loss, accounts for unpredictability across sample paths. Our policies trade off these two losses to get strong performance guarantees.

Our framework allows flexibility in choosing benchmarks. To understand why this is important, consider two common benchmarks for dynamic pricing: a controller has inventory B, and posts prices for T sequential customers, each of which has a random valuation. One common benchmark, known as the offline or prophet benchmark, considers a controller with full information of all randomness. It is easy to show that no online policy can get better than  $\Omega(T)$ regret against this benchmark. An alternate benchmark, known as the ex ante or fluid benchmark, corresponds to replacing all random quantities with their expectations. Here again, no online policy can get better than  $\Omega(\sqrt{T})$  regret (Vera and Banerjee 2020). Our approach, however, lets us identify benchmarks that have O(1) regret for all our settings.

Prophet and fluid benchmarks are also widely used in adversarial models of online allocation, leading to algorithms with worst-case guarantees. In contrast, we consider stochastic inputs, and consequently get much stronger guarantees. In particular, all our guarantees are parametric and depend explicitly on the distributions and problem primitives (i.e., constant parameters defining the instance). All our policies, however, have regret that is independent of the horizon and budgets.

# 2. Preliminaries and Overview2.1. Problem Settings and Results

We illustrate our framework by developing lowregret algorithms for the following problems:

Online stochastic knapsack. This serves as a baseline for our other problems. The controller has an initial resource budget B, and items arrive sequentially over T periods. Each item has a random type j,

which corresponds to a known resource requirement (or weight)  $w_j$  and a random reward  $R_j$ . In period t = T, T - 1, ..., 1 (where t denotes the time-to-go), we assume the arriving type is drawn from a finite set [n] from some known distribution  $\mathbf{p} = (p_1, ..., p_n)$ . At the start of each period, the controller observes the type of the arriving item, and must decide to accept or reject the item. The expected reward from selecting a type-j item is  $r_j = \mathbb{E}[R_j]$ .

Online probing. As before, an arriving type j has known expected reward  $r_j$ , but unknown realized reward  $R_j$ . Now the controller has the additional option of probing each request to observe the realization, and then accept/reject the item based on the revealed reward. The controller can also choose to accept the item without probing. In addition to the resource budget  $B_j$ , the controller has an additional probing budget  $B_j$  that limits the number of arrivals that can be probed. This introduces a trade-off between depleting the resource budget B and probing budget  $B_p$ . We assume here that  $R_j$  has finite support  $\{r_{jk}\}_{k \in [m]}$  of size m, and define  $q_{jk} := \mathbb{P}[R_j = r_{jk}]$  for  $k \in [m]$ . Note this reduces to OnlineKnapsack when either  $B_p \geq T$  or  $B_p = 0$ .

Dynamic pricing. The controller has an initial inventory  $B \in \mathbb{N}^d$  for d different resources. There are n types of customers, where a customer of type j requests a specific subset  $A_j \in \{0,1\}^d$  of resources, and has private valuation  $R^t \sim F_j$ . In each period t, the controller observes the customer type  $j \in [n]$ , and if sufficient resources are available, posts a price (fare) f from a finite set  $\{f_{j1}, \ldots, f_{jm}\}$ . The customer then purchases iff  $R^t > f$ . The vectors  $A_j$  and valuation functions  $(F_j : j \in [n])$  are known, but otherwise arbitrary. More generally, our technique handles probabilistic customer-choice models, where a customer, when presented with a price menu over bundles, picks a random bundle via some known distribution (which may depend on the menu).

Knapsack with distribution learning. We return to the OnlineKnapsack setting where items of type  $j \in [n]$  have weight  $w_j$  and random reward  $R_j$ . Now, however, the controller is unaware of the distribution of  $R_j$ , and must learn it from observations. In period t, the controller observes the arrival type j, and decides to accept/reject based on observed rewards up to time t. We consider two feedback models: full feedback, where the controller observes  $R_j$  regardless of whether the item is accepted or rejected, and censored feedback, where the controller only observes rewards of accepted items. For the latter (which is sometimes referred to as online contextual bandits with knapsacks), we assume the rewards  $R_j$  have sub-Gaussian tails (Boucheron et al. 2013, section 2.3).

Benchmarks and guarantees. Our framework, RABBI (for resolve and act based on the Bellman inequalities; see Section 3.2) is based on comparing two

controllers: Offline, which acts optimally given future information; and a nonanticipative controller Online, which tries to follow Offline. Both start in the same initial state  $S^T$ . We denote  $v^{\rm off}$  as the expected total reward collected by Offline acting optimally (i.e., according to a Bellman equation) given its information structure. In contrast, Online uses a nonanticipative policy  $\pi$  that maps current states to actions, resulting in a total expected reward  $v^{\rm on}_\pi$ .

Let  $\pi_R$  denote the online policy produced by our RABBI framework, and  $\pi$  denote any nonanticipative policy. Then the expected regret of  $\pi_R$  relative to the chosen offline benchmark is

$$\mathbb{E}[\mathsf{Regret}] \coloneqq v^{\mathsf{off}} - v^{\mathsf{on}}_{\pi_R} \ge \max_{\pi} [v^{\mathsf{on}}_{\pi}] - v^{\mathsf{on}}_{\pi_R}.$$

The last inequality, which follows from the fact that  $v_{\pi}^{\rm on} \leq v^{\rm off}$  for any pair of benchmark and online policies, emphasizes that the regret is a bound on the additive gap with respect to (w.r.t.) the best online policy.

For all of these problems, we use the RABBI framework to identify an appropriate benchmark with respect to which we get the following guarantees. First, for the OnlineKnapsack, we recover a result proved in Arlotto and Gurvich (2019) and Vera and Banerjee (2020).

**Theorem 1** (Theorem 1 in Arlotto and Gurvich 2019). For known reward distributions with finite mean, an online policy based on the RABBI framework obtains regret that depends only on the primitives  $(n, \mathbf{p}, \mathbf{r}, \mathbf{w})$ , but is independent of the horizon length T and resource budget B.

Theorem 1 builds intuition for using RABBI in more complex settings. In particular, the benchmark used in Theorem 1 is the full-information prophet, which is too loose for obtaining constant regret in the remaining settings (pricing, probing, and bandits; see Example 1). This is where our framework helps in guiding the choice of the right benchmark. In particular, we obtain the following results.

**Theorem 2** (Online Probing). For reward distributions with finite support of size m, an online-probing policy based on the rabbi framework (Algorithm 2) obtains regret that depends only on  $(n, m, \mathbf{q}, \mathbf{p}, \mathbf{r})$ , but is independent of horizon length T, resource budget B and probing budget  $B_p$ .

**Theorem 3** (Dynamic Pricing). For any reward distributions  $(F_j : j \in [n])$  and prices  $\mathbf{f}$ , a pricing policy based on the RABBI framework (Algorithm 3) obtains regret that depends only on  $(A, \mathbf{f}, F_1, \ldots, F_n)$ , but is independent of horizon length T and initial budget levels  $B \in \mathbb{N}^d$ .

The result for dynamic pricing also extends naturally to resource bundles and general customer-choice models (see Section 5.5 and Theorem 6 therein).

For the bandit settings, we define a separation parameter  $\delta = \min_{j \neq j'} \mathbb{E}[R_j]/w_j - \mathbb{E}[R_{j'}]/w_{j'}$ ; this is only for our bounds, and is not known to the algorithm.

**Theorem 4** (Knapsack with Distribution Learning). Assuming the reward distributions are sub-Gaussian, in the full feedback setting, a policy based on the RABBI framework (Algorithm 5) obtains regret that depends only on the primitives  $(n, \mathbf{p}, \mathbf{r}, \mathbf{w}, \delta)$  and is independent of the horizon length T and knapsack capacity B.

The last result can also be used as a black box for the censored feedback setting to get an  $O(\log T)$  regret guarantee (see Corollary 1 in Section 6.3).

### 2.2. Overview of Our Framework

We develop our framework in the full generality of MDPs in Section 3. To give an overview and gain insight into the general version, we use OnlineKnapsack as a warm-up. A schema for the framework is provided in Figure 1.

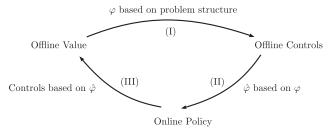
In the OnlineKnapsack problem, at any time-to-go t, let  $Z_j^t \in \mathbb{N}$  denote the (random) number of type-j arrivals in the remaining t periods. Recall rewards of type j arrivals have expected value  $r_j := \mathbb{E}[R_j]$ . Define Offline to be a controller that knows  $Z^t$  for all t in advance. The total reward collected by Offline can be written as an integer linear program:

$$V(t, b|Z^{t}) = \max_{\mathbf{x}_{a} \in \mathbb{N}^{n}} \{ \mathbf{r}' \mathbf{x} : \mathbf{w}' \mathbf{x}_{a} \le b, \mathbf{x}_{a} \le Z^{t} \}$$
$$= \max_{\mathbf{x}_{a}, \mathbf{x}_{r} \in \mathbb{N}^{n}} \{ \mathbf{r}' \mathbf{x}_{a} : \mathbf{w}' \mathbf{x}_{a} \le b, \mathbf{x}_{a} + \mathbf{x}_{r} = Z^{t} \}. \quad (1)$$

The function  $V(\cdot|Z^t)$  is thus Offline's value function (Figure 1), where the notation  $|Z^t|$  emphasizes that V is conditioned on  $Z^t$ . Moreover, for every j, the variables  $x_{a,j}, x_{r,j}$  represent action summaries: the number of type-j arrivals accepted and rejected, respectively.

The value  $V(\cdot|Z^t)$  can also be represented via Bellman equations. Specifically, at time-to-go t, assuming

Figure 1. The RABBI Framework



*Notes.* We first define Offline's value function by specifying access to future information. Next, we identify a tractable relaxation  $\varphi$  for Offline's value under this same information structure (step I). Finally, we introduce a nonanticipative estimate  $\hat{\varphi}$  for  $\varphi$ , and use it to design online controls (step II). The resulting online policy is evaluated against Offline's value (step III).

Offline has budget b and the arriving type is  $\xi$ , the value function obeys the Bellman equation:

$$V(t,b|Z^{t}) = \max \{ [r_{\xi^{t}} + V(t-1,b-w_{\xi^{t}}|Z^{t-1})] \mathbb{1} w_{\xi^{t}} \le b,$$

$$V(t-1,b|Z^{t-1}) \}, \quad \forall t,b,\xi^{t}.$$

Next, consider the linear programming relaxation for V(t, b):

$$\varphi(t, b|Z^t) := \max_{\mathbf{x}_a, \mathbf{x}_r > 0} \{ \mathbf{r}' \mathbf{x}_a : \mathbf{w}' \mathbf{x}_a \le b, \mathbf{x}_a + \mathbf{x}_r = Z^t \}.$$

It is clear that  $\varphi$  is more tractable than V, and that it approximates V up to an integrality gap. However,  $\varphi$  does not obey a Bellman equation. To circumvent this, we introduce the notion of Bellman inequalities, wherein we require that  $\varphi$  satisfies Bellman-like conditions for most sample paths. Formally, for some random variables  $L_B$ , we want  $\varphi$  to satisfy

$$\varphi(t, b|Z^{t}) \leq \max \{ [r_{\xi^{t}} + \varphi(t-1, b-w_{\xi^{t}}|Z^{t-1})] \mathbb{1}_{\{w_{\xi^{t}} \leq b\}}, \\ \varphi(t-1, b|Z^{t-1}) \} + L_{B}(t, b).$$

Note that, if  $\mathbb{E}[L_B(t,b)]$  is small, with expectation taken over  $Z^t$ , then  $\varphi$  almost satisfies the Bellman equations. We henceforth refer to  $\varphi$  as a relaxed value for V and  $L_B$  as the Bellman loss.

Establishing that actions derived from  $\varphi$  are nearly optimal for Offline accomplishes step (I) in Figure 1. For step (II), we want to emulate Offline by estimating  $\varphi$  based on current information. A natural estimate is obtained by taking expectations over future randomness, to get the following:

$$\hat{\varphi}(t,b) := \max_{\mathbf{y}_{\mathtt{a}},\mathbf{y}_{\mathtt{r}} \geq 0} \Big\{ \mathbf{r}'\mathbf{y}_{\mathtt{a}} : \mathbf{w}'\mathbf{y}_{\mathtt{a}} \leq b, \mathbf{y}_{\mathtt{a}} + \mathbf{y}_{\mathtt{r}} = \mathbb{E}\big[Z^t\big] \Big\}.$$

Note that  $\hat{\varphi}$  does not approximate V or  $\varphi$  up to a constant additive error (Vera and Banerjee 2020); however,  $\hat{\varphi}$  can be used as a predictor for the action taken by Offline. Specifically, at time t with current budget b, rabbi first computes  $\hat{\varphi}(t,b)$  and then interprets the solution  $\mathbf{y}$  as a score for each action (here, accept/reject). We show that taking the action with the highest score (i.e., action  $\arg\max_{u\in\{\mathbf{a},\mathbf{x}\}}\{y_{\xi^t,u}\}$ ) guarantees that Online and Offline play the same action with high probability. Whenever Offline and Online play different actions, we incur a loss, which we refer to as the information loss, as it quantifies how having less information impacts Online's actions. This process of using  $\hat{\varphi}$  to derive actions is represented as step (III) in Figure 1.

**2.2.1. Toward a General Framework.** For all the problems in Section 2.1, our approach uses a similar three-step process, wherein we choose an Offline benchmark,

identify relaxed value  $\varphi$  via appropriate optimization problem, and get an online policy based on estimate  $\hat{\varphi}$ . Consequently, we refer to our framework as RABBI, which stands for resolve and act based on Bellman inequalities.

Our work builds on constant-regret policies for multidimensional packing (Vera and Banerjee 2020), and more general online optimization problems (Banerjee and Freund 2020). The techniques developed in these works, however, have two fundamental shortcomings that prevent them from addressing the settings we consider:

- Use of full-information benchmarks: Existing works (Arlotto and Gurvich 2019, Banerjee and Freund 2020, Vera and Banerjee 2020) use the full-information benchmark, which is too loose for our settings. Indeed, for probing/pricing/learning settings, no algorithm can have constant regret compared with the full-information benchmark (see Example 1).
- Explicit value function characterizations: The optimization problem in Equation (1) has a closed-form solution, which was used explicitly by Arlotto and Gurvich (2019), Vera and Banerjee (2020), and Banerjee and Freund (2020). This does not extend to more complex settings.

Our framework in this work resolves these short-comings in a structured way, allowing us to get provably near-optimal algorithms for several canonical resource-allocation problems. Moreover, we do so via a generalized notion of information-augmented benchmarks, and our decomposition of the regret into the information loss (capturing randomness in inputs) and Bellman loss (capturing limited computational power). This flexibility helps greatly in the design of our algorithms.

#### 2.3. Related Work

Our approach has commonalities with two closely related approaches:

Prophet inequalities and ex ante relaxations: A well-studied framework for obtaining performance guarantees for heuristics policies is to compare against a full-information agent, or prophet. This line of work focuses on competitive-ratio bounds (see Kleinberg and Weinberg 2012, Correa et al. 2017, Düetting et al. 2017 for overviews of the area). In particular, Correa et al. (2017) obtains a multiplicative guarantee for dynamic posted pricing with a single item under worst-case distribution. A related line of work considers the use of ex ante linear program (LP) relaxations (Alaei 2014, Buchbinder et al. 2014) for obtaining worst-case competitive guarantees in online packing problems. In contrast, we obtain an additive guarantee for multiple items in a parametric setting.

MDP dual relaxations: A standard way to get bounds on MDPs is via information-relaxations, which at a

high level, create benchmarks by endowing Offline with additional information, while forcing it to pay a penalty for using this information. Brown et al. (2010) and Balseiro and Brown (2019) use this in a dual-fitting approach to construct performance bounds for greedy algorithms in different problems. In contrast, our framework is similar to a primal-dual approach: we adaptively construct our relaxations, and derive controls directly from them. We compare the two approaches in more detail in Online Appendix EC.5. Moreover, the different problems we apply RABBI to each have a large body of prior work.

**2.3.1. Online Packing.** There is a long line of work on the baseline OnlineKnapsack and generalizations. A notable work in this line is Jasin and Kumar (2012), who gives a policy with constant expected regret when the problem instance is far from a set of certain nondegenerate instances. This inefficiency, though, is fundamental, since they use the ex ante (or fluid) benchmark, which has  $\Omega(\sqrt{T})$  under nondegeneracy. More recently, Bumpensanti and Wang (2020) partially extend the result of Arlotto and Gurvich (2019) for more general packing problems; however, their policy only gives constant regret under independent and identically distributed (i.i.d.) Poisson arrivals, and requires the system to be scaled linearly (i.e., *B* grows proportional to *T*). In contrast, Arlotto and Gurvich (2019) (one dimension) and Vera and Banerjee (2020) (multiple dimensions) provide constant-regret policies with no assumption on the scaling. The approach in the latter is further generalized in Banerjee and Freund (2020) to handle more complex problems including bin packing and quality of service (QOS) constraints. See Vera and Banerjee (2020), Banerjee and Freund (2020) for more discussion and references.

**2.3.2. Probing.** Approximation algorithms have been developed for offline probing problems, both under budget constraints (Gupta and Nagarajan 2013) and probing costs (Weitzman 1979, Singla 2018). Another line of work pursues tractable nonadaptive constant-factor competitive algorithms for this problem (Gupta et al. 2016). In terms of online adaptive algorithms, Chugg and Maehara (2019) introduce an algorithm with bounded competitive ratio in an adversarial setting.

**2.3.3. Dynamic Posted Pricing.** This is a canonical problem in operations management with a vast literature (see Talluri and Van Ryzin 2006 for an overview). Much of this literature focuses on asymptotically optimal policies in regimes where the inventory B and/or horizon T grow large. When B and T are scaled together by a factor k, there are known algorithms with regret that scales as  $O(\sqrt{k})$  or  $O(\log(k))$ , depending on assumptions on the primitives (e.g., smoothness of the

demand with price) (Jasin 2014). There is also vast literature on pricing when the demand function is not known and has to be learned (Chen et al. 2019). Finally, under adversarial arrivals, Babaioff et al. (2015) provide a policy with  $O((B \log T)^{2/3})$  regret under adversarial inputs, as opposed to our O(1) guarantee under stochastic inputs.

**2.3.4. Knapsack with Learning.** Multiarmed bandit problems have been widely studied, and we refer to Bubeck et al. (2012, 2013) for an overview. Bandit problems with combinatorial constraints on the arms are known as *bandits with knapsacks* (Badanidiyuru et al. 2018), and the generalization where arms arrive online is known as *contextual bandits with knapsacks* (Badanidiyuru et al. 2014, Agrawal and Devanur 2016). Results in this literature typically study worst-case distributions. We, in contrast, pursue parametric regret bounds that explicitly depend on the (unknown) discrete distribution. Closest to our work is Wu et al. (2015), who provide a upper confidence bounds (UCB)-based algorithm that gets  $O(\sqrt{T})$  regret (in contrast, we get  $O(\log T)$  regret for the same setting).

# 3. Approximate Control Policies via the Bellman Inequalities

In this section, we describe our general framework. Before proceeding, we introduce some notation. We work an underlying probability space  $(\Omega, \Sigma, \mathbb{P})$ , and for any event  $\mathcal{B} \subseteq \Omega$ , we denote its complement by  $\mathcal{B}^c$ . We use boldface letters to indicate vector-valued variables (e.g.  $\mathbf{p}, \mathbf{w}$ , etc.), and capital letters to denote matrices and/or random variables. For an optimization problem (P), we use P to denote its optimal value. When using LP formulations with decision variables  $\mathbf{x}$ , we interchangeably use  $x_{ij} = x(i,j)$  to denote the (i,j)th component of  $\mathbf{x}$ .

### 3.1. Offline Benchmarks and Bellman Inequalities

We consider an online decision-making problem with state space  $\mathcal S$  and action space  $\mathcal U$ , evolving over periods  $t=T,T-1,\ldots,1$ ; here T denotes the horizon, and t is the time to-go. In any period t, the controller first observes a random arrival  $\xi^t \in \Xi$ , following which it must choose an action  $u \in \mathcal U$ . For system-state  $s \in \mathcal S$  at the beginning of period t, and random arrival  $\xi \in \Xi$ , an action  $u \in \mathcal U$  results in a reward  $\mathcal R(s,\xi,u)$ , and transition to the next state  $\mathcal T(s,\xi,u)$ . We assume both reward and future state are random variables whose realizations are determined for every u given  $\xi$ . This assumption is for ease of exposition only; our results can be extended to hold when rewards or transitions are random given  $\xi$ .

The feasible actions for state s and input  $\xi$  correspond to the set  $\{u \in \mathcal{U} : \mathcal{R}(s, \xi, u) > -\infty\}$ . We assume that this feasible set is nonempty for all  $s \in \mathcal{S}, \xi \in \Xi$ ,

and that the maximum reward is bounded, that is,  $\sup_{s \in \mathcal{S}, \xi \in \Xi, u \in \mathcal{U}} \mathcal{R}(s, \xi, u) < \infty$ .

The MDP described previously induces a natural filtration  $\mathcal{F}$ , with  $\mathcal{F}_t = \sigma(\{\xi^\tau : \tau \ge t\})$ ; a nonanticipative policy is one that is adapted to  $\mathcal{F}_t$ . We allow Offline to use a richer information filtration  $\mathcal{G}$ , where  $\mathcal{G}_t \supseteq \mathcal{F}_t$ . Note that since t denotes the time-to-go, we have  $\mathcal{G}_{t-1} \supseteq \mathcal{G}_t$ . Henceforth, to keep track of the information structure, we use the notation  $f(\cdot|\mathcal{G}_t)$  to clarify that a function f is measurable with respect to the sigma field  $\mathcal{G}_t$ .

Given any filtration  $\mathcal{G}$ , Offline is assumed to play the optimal policy adapted to  $\mathcal{G}$ , hence Offline's value function is given by the following Bellman equation:

$$V(t, s|\mathcal{G}_t) = \max_{u \in \mathcal{U}} \{ \mathcal{R}(s, \xi^t, u) + \mathbb{E}[V(t-1, \mathcal{T}(s, \xi^t, u)|\mathcal{G}_{t-1})|\mathcal{G}_t] \}, \quad (2)$$

with the boundary condition  $V(0,\cdot) = 0$ . We denote the expected value as  $v^{\text{off}} := \mathbb{E}[V(T,S^T|\mathcal{G}_T)]$ . Note that  $v^{\text{off}}$  is an upper bound on the performance of the optimal nonanticipative policy.

We present a specific class of filtration (generated by augmenting the canonical filtration) that suffice for our applications (see Figure 2 for an illustration of the definition).

**Definition 1** (Canonical Augmented Filtration). Let  $G_{\Theta} := (G_{\theta} : \theta \in \Theta)$  be a set of random variables. The canonical filtration w.r.t.  $G_{\Theta}$  is

$$\mathcal{G}_t = \sigma(\{\xi^l : l \geq t\} \cup G_{\Theta}) \supseteq \mathcal{F}_t.$$

The richest augmented filtration is the full-information filtration, wherein  $\mathcal{G}_t = \mathcal{F}_1$  for all t, that is, the canonical filtration with  $G_{\Theta} = (\xi^t : t \in [T])$ . As  $\mathcal{G}_t$  gets coarser, the difference in performance between Offline and Online decreases. Indeed, when  $\mathcal{G} = \mathcal{F}$ , then Equation (2) reduces to the Bellman equation for the value function of an optimal nonanticipative policy:

$$\begin{split} V(t,s|\mathcal{F}_t) &= \max_{u \in \mathcal{U}} \left\{ \mathcal{R}\left(s,\xi^t,u\right) \right. \\ &+ \mathbb{E}\left[V\left(t-1,\mathcal{T}\left(s,\xi^t,u\right)|\mathcal{F}_{t-1}\right)\right] \right\}, \quad V(0,\cdot,\cdot) = 0, \end{split}$$

where the expectation is taken with respect to the next period's input  $\xi^{t-1}$ .

**Example 1** (Full Information Is Too Loose). Consider a dynamic pricing instance with n=d=1, prices  $\mathbf{f}=(1,2)$ , and valuation distribution  $\mathbb{P}[R^t=1+\varepsilon]=p$  and  $\mathbb{P}[R^t=2+\varepsilon]=1-p$ . When B=T, the optimal policy always posts a price that maximizes  $(f\cdot\mathbb{P}[R^t>f])$ . If  $p\geq 1/2$ , then the optimal policy (DP) always posts price f=1 and has expected reward T. On the other hand, full information can post price  $R^t-\varepsilon$  at time t and extract full surplus  $v^{\text{off}}=\sum_t\mathbb{E}[R^t-\varepsilon]=T(2-p)$ . Thus, the regret against full information must grow as  $\Omega(T)$ . This example is not pathological; the same behavior persists even in random instances (see Section 5.4).

We are now ready to introduce the notion of relaxed value  $\varphi$  and Bellman inequalities. Intuitively,  $\varphi$  is almost defined by a dynamic-programming recursion; quantitatively, whenever  $\varphi$  does not satisfy the Bellman equation, we incur an additional loss  $L_B$ , which we denote the Bellman loss.

**Definition 2** (Bellman Inequalities). The family of random variables  $\{\varphi(t,s)\}_{t,s}$  satisfies the Bellman inequalities w.r.t. filtration  $\mathcal{G}$  and random variables  $\{L_B(t,s)\}_{t,s}$  if  $\varphi(t,\cdot)$  and  $L_B(t,\cdot)$  are  $\mathcal{G}_t$ -measurable for all t and the following conditions hold:

- 1. Initial ordering:  $\mathbb{E}[V(T, S^T)|\mathcal{G}_T] \leq \varphi(T, S^T|\mathcal{G}_T)$ .
- 2. Monotonicity:  $\forall s \in \mathcal{S}, t \in [T]$ ,

$$\varphi(t, s|\mathcal{G}_t) \leq \max_{u \in \mathcal{U}} \{ \mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u)|\mathcal{G}_{t-1})|\mathcal{G}_t] \} + L_B(t, s).$$
(3)

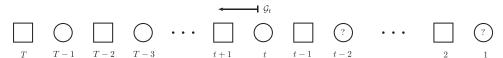
3. Terminal condition:  $\varphi(0,s) = 0 \ \forall s \in \mathcal{S}$ .

We refer to  $\varphi$  and  $L_B$  as the relaxed value and Bellman loss pair with respect to  $\mathcal{G}$ , and use  $|\mathcal{G}_t|$  to remind the reader that we need the information contained in  $\mathcal{G}_t$  to evaluate  $\varphi(t,s)$ .

Given any  $\varphi$ , monotonicity holds trivially with  $L_B = \varphi$  (but leads to poor performance guarantees). On the other hand,  $\varphi$  (which may be intractable) is the only value function guaranteeing  $L_B = 0$ . The crux of our approach is to identify a good  $\varphi$  balances the loss and tractability.

A special case is when the Bellman loss is 0 over sample paths in some chosen set.

Figure 2. Illustration of Definition 1



*Notes.* In online probing (see Section 4), arrivals first reveal their public type, then the controller chooses an action (accept/probe/reject), and then the private type (true reward) is revealed. Squares (respectively, circles) represent public (respectively, private) information. The filtration  $\mathcal{G}$  used by RABBI comprises of all public types, that is,  $G_{\Theta} = (\xi^{\theta} : \xi^{\theta} \text{ is a public type})$ . At time t, Offline knows all the information thus far (to the left and including t), plus the future squares.

**Definition 3** (Exclusion Sets). A set  $\mathcal{B}(t,s)$  is an exclusion set if we can write the Bellman loss as  $L_B(t,s) = r_{\varphi} \mathbb{1}_{\mathcal{B}(t,s)}$  for some constant  $r_{\varphi} > 0$  and events  $\mathcal{B}(t,s) \subseteq \Omega$ .

If the Bellman loss can be defined with exclusion sets, then from Definition 2 (monotonicity), we obtain the condition  $\varphi(t,s|\mathcal{G}_t) \leq \max_{u \in \mathcal{U}} \{\mathcal{R}(s,\xi^t,u) + \mathbb{E}[\varphi\} (t-1,\mathcal{T}(s,\xi^t,u)|\mathcal{G}_{t-1})|\mathcal{G}_t]$ , that is, monotonicity is satisfied for all realizations  $\omega \in \Omega$  except for those in the exclusion set  $\mathcal{B}(t,s)$ .

To build intuition, we specify the Bellman inequalities for our baseline OnlineKnapsack. For this end, we first need the following lemma characterizing the sensitivity of LP solutions.

**Lemma 1.** Consider an LP  $(P[\mathbf{d}])$ :  $\max\{\mathbf{r'x} : M\mathbf{x} = \mathbf{d}, \mathbf{x} \ge 0\}$ , where  $M \in \mathbb{R}^{m \times n}$  is an arbitrary constraint matrix. If  $\bar{\mathbf{x}}$  solves  $(P[\mathbf{d}])$  and  $\bar{x}_j \ge 1$  for some j, then  $P[\mathbf{d}] = r_j + P[\mathbf{d} - M_j]$ .

**Proof.** By assumption, the optimal value of  $(P[\mathbf{d}])$  remains unchanged if we add the inequality  $x_j \ge 1$ . Therefore, we have  $P[\mathbf{d}] = \max\{\mathbf{r}'(\mathbf{x} + \mathbf{e}_i) : M(\mathbf{x} + \mathbf{e}_i) = \mathbf{d}, \mathbf{x} \ge 0\}$ .  $\square$ 

Lemma 1 lets us divide  $P[\mathbf{d}]$  into two summands: the immediate reward  $r_j$  and the future reward  $P[\mathbf{d} - M_j]$ . This has the flavor of dynamic programming we need for defining the Bellman loss.

**Example 2** (Bellman Loss for Baseline Setting). For the baseline OnlineKnapsack, discussed in Section 2.2, we chose the full-information filtration  $\mathcal{G}_t = \mathcal{F}_1$  for all t so that  $\varphi(t,b|\mathcal{G}_t) := \max_{\mathbf{x} \geq 0} \{\mathbf{r'}\mathbf{x}_a : \mathbf{w'}\mathbf{x}_a \leq b, \mathbf{x}_a + \mathbf{x}_r = Z^t\}$ . We define the exclusion sets as

$$\mathcal{B}(t,b) = \{ \omega \in \Omega : \exists \mathbf{x} \text{ solving } \varphi(t,b) \text{ s.t. } x(\mathbf{a},\xi^t) \ge 1$$
 or  $x(\mathbf{r},\xi^t) \ge 1 \}.$ 

By Lemma 1, outside the exclusion sets  $\mathcal{B}(t,b)$ , monotonicity holds with zero Bellman Loss, that is,

$$\varphi(t, s|\mathcal{G}_t) \leq \max_{u \in \mathcal{U}} \{ \mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u)|\mathcal{G}_{t-1})|\mathcal{G}_t] \} \quad \forall \omega \notin \mathcal{B}(t, s).$$

Moreover, for our choice of  $\varphi$ , since the optimal solution sorts items by  $r_j/w_j$ , we have that the maximum loss outside the exclusion set is bounded by  $r_{\varphi} \leq \max_{j,i} \{w_i r_j/w_j - r_i\}$ , which depends only on the primitives. Thus, Definition 2 is satisfied with Bellman loss  $L_B(t,b) = r_{\varphi} \mathbb{1}_{B(t,b)}$ .

To generalize this, we need two definitions. First, we define the maximum Bellman loss as follows.

**Definition 4** (Maximum Loss). For a given relaxation  $\varphi$ , the maximum loss is given by

$$\begin{split} r_{\varphi} &:= \max_{t,s,u:\mathcal{R}(s,\xi^t,u)>-\infty} \big\{ \varphi(t,s|\mathcal{G}_t) - \big(\mathcal{R}\big(s,\xi^t,u\big) \\ &+ \mathbb{E}\big[ \varphi\big(t-1,\mathcal{T}\big(s,\xi^t,u\big)|\mathcal{G}_{t-1}\big)|\mathcal{G}_t \big] \big) \big\}. \end{split}$$

Next, note that the optimal action in the right-hand side (RHS) of Equation (3) need not be unique, and indeed the inequality can be satisfied by multiple actions. For given  $\varphi$  and  $L_B$ , we define the following.

**Definition 5** (Satisfying Actions). Given a filtration  $\mathcal{G}$  and relaxed value  $\varphi$ , we say that u is a satisfying action for state s at time t if

$$\varphi(t, s|\mathcal{G}_t) \le \mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u)|\mathcal{G}_{t-1})|\mathcal{G}_t] + L_B(t, s).$$
(4)

At any time t and state  $s \in \mathcal{S}$ , any action in  $\arg\max_{u \in \mathcal{U}} \{\mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u)|\mathcal{G}_{t-1})|\mathcal{G}_t]\}$  is always a satisfying action (see monotonicity in Definition 2). Moreover, to identify a satisfying action, we must know  $\mathcal{G}_t$ . We now have the following proposition.

**Proposition 1.** Consider a relaxation  $\varphi$  and Bellman loss  $L_B$  that satisfy the Bellman inequalities w.r.t. filtration  $\mathcal{G}$ . Let  $(S^t, t \in [T])$  denote the state trajectory under a policy that, at time t, takes any satisfying action  $U^t = U^t(S^t|\mathcal{G}_t)$ . Then,

$$\mathbb{E}[V(T,S^T|\mathcal{G}_T)] - \mathbb{E}\left[\sum_{t=1}^T \mathcal{R}(S^t,\xi^t,U^t)\right] \leq \mathbb{E}\left[\sum_{t=1}^T L_B(t,S^t|\mathcal{G}_t)\right].$$

**Proof.** From the monotonicity condition in the Bellman inequalities (Definition 2), and the definition of a satisfying action (Definition 5), we have, for all time t, that

$$\varphi(t, S^t | \mathcal{G}_t) \leq \mathbb{E} [\mathcal{R}(S^t, \xi^t, U^t) + \varphi(t - 1, S^{t-1} | \mathcal{G}_{t-1}) + L_B(t, S^t | \mathcal{G}_t) | \mathcal{G}_t].$$

Iterating this inequality over t, we get  $\varphi(T, S^T | \mathcal{G}_T) \leq \sum_{t=1}^T \mathbb{E}[\mathcal{R}(S^t, \xi^t, U^t) + L_B(t, S^t | \mathcal{G}_t) | \mathcal{G}_t]$ . Finally, by the initial ordering condition, we have  $\mathbb{E}[V(T, S^T) | \mathcal{G}_T] \leq \varphi(T, S^T | \mathcal{G}_T)$ .  $\square$ 

Proposition 1 shows that a policy that always plays a satisfying action  $U^t$  approximates the performance of Offline up to an additive gap given by the total Bellman loss  $\mathbb{E}\left[\sum_{t=1}^T L_B(t, S^t | \mathcal{G}_t)\right]$ . More importantly, it suggests that Online should try to track Offline by guessing and playing a satisfying action  $U^t$  in each period. We next illustrate how Online can generate such guesses.

#### 3.2. From Relaxations to Online Policies

Suppose we are given an augmented canonical filtration  $\mathcal{G}_t = \sigma(\{\xi^I : l \geq t\} \cup G_\Theta)$ , and assume that the relaxed value  $\varphi$  can be represented as a function of the random variables  $\{\xi^I : l \geq t\} \cup G_\Theta$  as  $\varphi(t, s | \mathcal{G}_t) = \varphi(t, s; f_t(\xi^T, \dots, \xi^t, G_\Theta))$ . In particular, we henceforth focus on a special case where  $\varphi$  is expressed as the solution of an optimization problem:

$$\varphi(t,s;f_t(\xi^T,\ldots,\xi^t,G_{\Theta})) = \max_{\mathbf{x}\in\mathbb{R}^{t\ell\times\mathbb{Z}}} \left\{ h_t(\mathbf{x};s,f_t(\xi^T,\ldots,\xi^t,G_{\Theta})) : g_t(\mathbf{x};s,f_t(\xi^T,\ldots,\xi^t,G_{\Theta})) \le 0 \right\}$$
(5)

The decision variables give action summaries: for given state s and time t,  $x_{u,\xi}$  represents the number of times action u is taken for input  $\xi$  in remaining periods. We can also interpret  $x_{u,\xi}$  as a score for action u when input  $\xi$  is presented. Now, to get a nonanticipative policy, a natural projection of  $\varphi(t,s|\mathcal{G}_t)$  on the filtration  $\mathcal{F}$  is given via the following optimization problem:

$$\hat{\varphi}(t,s|\mathcal{F}_t) = \varphi(t,s;\mathbb{E}[f_t(\xi^T,\ldots,\xi^t,G_{\Theta})|\mathcal{F}_t])$$

$$= \max_{\mathbf{y} \in \mathbb{R}^{U \times \mathbb{E}}} \{h_t(\mathbf{y};s,\mathbb{E}[f_t|\mathcal{F}_t]) : g_t(\mathbf{y};s,\mathbb{E}[f_t|\mathcal{F}_t]) \le 0\}.$$
(6)

The solution of this optimization problem gives action summaries (or scores) y. The main idea of the RABBI algorithm is to play the action with the highest score.

**Algorithm 1 RABBI** (Resolve and Act Based on Bellman Inequalities)

**Input:** Access to functions  $f_t$  such that  $\varphi(t,s|\mathcal{G}_t) = \varphi(t,s;f_t(\xi^T,\ldots,\xi^t,G_{\Theta})).$ 

**Output:** Sequence of decisions  $\hat{U}^t$  for Online.

- 1: Set  $S^T$  as the given initial state
- 2: **for** t = T, ..., 1 **do**
- 3: Compute  $\hat{\varphi}(t, S^t) = \varphi(t, S^t; \mathbb{E}[f_t(\xi^T, \dots, \xi^t, G_{\Theta}) | \mathcal{F}_t]) \text{ with associated scores } \mathbf{y} = \{y_{u,\xi}\}_{u \in \mathcal{U}, \xi \in \Xi}.$
- 4: Given input  $\xi^t$ , choose the action  $\hat{U}^t$  with the highest score  $y_{u,\xi^t}$ .
- 5: Collect reward  $\mathcal{R}(S^t, \xi^t, \hat{U}^t)$ ; update state  $S^{t-1} \leftarrow \mathcal{T}(S^t, \xi^t, \hat{U}^t)$ .

**Theorem 5.** Let Offline be defined by an augmented filtration  $\mathcal{G}_t$  as in Definition 1. Assume the relaxation  $\varphi(t,s)$  satisfies the Bellman inequalities with loss  $L_B$ , and for all  $(t,a) \in [T] \times \mathcal{S}$ , let  $\mathcal{Q}(t,s) \subseteq \Omega$  denote the set of sample paths where the action  $\hat{U}^t$  taken by RABBI is not a satisfying action. If  $(S^t, t \in [T])$  denotes the state trajectory under RABBI, then

$$\mathbb{E}[\text{Regret}] \leq \mathbb{E}\left[\sum_{t} \left(r_{\varphi} \mathbb{1}_{\mathcal{Q}(t,S^{t})} + \mathbb{1}_{\mathcal{Q}(t,S^{t})^{c}} L_{B}(t,S^{t})\right)\right]$$

$$\leq \sum_{t} \left(r_{\varphi} \mathbb{P}[\mathcal{Q}(t,S^{t})] + \mathbb{E}[L_{B}(t,S^{t})]\right).$$

**Remark 1** (Bellman and Information Loss). The bound in Theorem 5 has two distinct summands: The information loss  $\Sigma_t \mathbb{P}[\mathcal{Q}(t,S^t)]$  measures how often RABBI takes a nonsatisfying action due to randomness in sample paths; and the Bellman loss  $\Sigma_t \mathbb{E}[L_B(t,S^t)]$ ) quantifies violations of the Bellman equations made under the pseudo value function  $\varphi$ .

The proof of Theorem 5 is based on the compensated coupling approach introduced in Vera and Banerjee (2020). The idea is to imagine simulating controllers Offline and Online with identical random inputs ( $\xi^t$ :  $t \in [T]$ ), with Online acting before Offline.

Moreover, suppose at some time *t*, both controllers are in the same state s. Recall that, for any given state s at time t, an action u is satisfying if Offline's value does not decrease when playing u (Definition 5). If Online chooses to play a satisfying action, then we can make Offline play the same action, and consequently both move to the same state. On the other hand, if Online chooses an action that is not satisfying, then the two trajectories may separate. We can avoid this, however, by compensating Offline so that it agrees to take the same action as Online. In particular, it is always sufficient to compensate Offline by the maximum loss  $r_{\varphi}$  to ensure its reward does not decrease by following Online. As a consequence, the (compensated) Offline and Online take the same actions, and thus their trajectories are coupled.

As an example, for OnlineKnapsack with budget B=2, weights  $w_j=1$   $\forall$  j, and horizon T=5, consider a sample path  $\omega \in \Omega$  with rewards  $(\xi^5, \xi^4, \xi^3, \xi^2, \xi^1)=(5,7,2,7,2)$ . The sample path comprises three different types, and the sequence of actions (r,a,r,a,r) (selecting the value 7 items) is optimal for Offline, with total reward of 14. Suppose Online, in period t=5, wants to accept the item with reward  $\xi^5=5$ ; then, Offline is willing to follow this action if given a compensation of 2 (in addition to collecting reward 5). Offline and Online then start the next period t=4 in the same state with budget 1, hence remain coupled.

**Proof of Theorem 5.** Denoting Offline's state as  $\bar{S}^t$ , we have via Proposition 1 that  $\forall t$ :

$$\varphi(t, \bar{S}^t | \mathcal{G}_t) \leq \mathbb{E} \left[ \mathcal{R}(\bar{S}^t, \xi^t, U^t) + \varphi(t - 1, \bar{S}^{t-1} | \mathcal{G}_{t-1}) + L_B(t, \bar{S}^t) | \mathcal{G}_t \right].$$

Let us assume as the induction hypothesis that  $\bar{S}^t = S^t$ . This holds for t = T by definition. At any time t and state  $S^t$ , if  $\hat{U}^t$  is not a satisfying action for Offline, then we have from the definition of the maximum loss (Definition 4) that:

$$r_{\varphi} \geq \varphi(t, S^{t}|\mathcal{G}_{t}) - \mathcal{R}(S^{t}, \xi^{t}, \hat{U}^{t}) + \mathbb{E}[\varphi(t-1, S^{t-1}|\mathcal{G}_{t-1})|\mathcal{G}_{t}]) \quad almost \ surely.$$

Now to make Offline take action  $\hat{U}^t$  so as to have the same subsequent state as Online, it is sufficient to compensate Offline with an additional reward of  $r_{\varphi}$ . Specifically, we have

$$\varphi(t, S^t | \mathcal{G}_t) \leq \mathbb{E} \left[ \mathcal{R} \left( S^t, \xi^t, \hat{\mathcal{U}}^t \right) + \varphi(t - 1, S^{t-1} | \mathcal{G}_{t-1}) \right. \\ + r_{\varphi} \mathbb{1}_{\mathcal{Q}(t, S^t)} + \mathbb{1}_{\mathcal{Q}(t, S^t)^c} L_B(t, S^t) | \mathcal{G}_t \right].$$

Finally, as in Proposition 1, we can iterate over t to obtain

$$\mathbb{E}[\varphi(T, S^T | \mathcal{G}_T)] \leq \mathbb{E}\left[\sum_t \mathcal{R}(S^t, \xi^t, \hat{\mathcal{U}}^t). + \sum_t \left(r_{\varphi} \mathbb{1}_{\mathcal{Q}(t, S^t)} + \mathbb{1}_{\mathcal{Q}(t, S^t)^c} L_B(t, S^t)\right)\right].$$

The first sum on the right-hand side corresponds exactly to Online's total reward using the rabbi policy. By the initial ordering property,  $\mathbb{E}[V(T,S^T)] \leq \mathbb{E}[\varphi(T,S^T)]$ , and we get the result.  $\square$ 

### 4. Online Probing

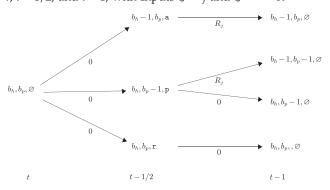
We now apply our framework to online probing. Here, each arrival type *j* has an independent random reward  $R_i \in \{r_{ik} : k \in [m]\}$  drawn with probabilities  $\{q_{ik}\}$ ; **r** and **q** are known. We assume without loss of generality (w.l.o.g.) that  $r_{i1} < r_{i2} < \ldots < r_{im}$  and  $r_{im} > 0$ . For ease of exposition, we assume that all arrivals have unit weights; our analysis however extends to general weights  $w_i$ . The controller may accept (a), reject (r), or probe (p) the arrival. Accepting a type-j item without probing results in expected reward of  $\bar{r}_i := \sum_{k \in [m]} r_{ik} q_{ik}$ . Probing reveals the realized reward, after which it can be accepted or rejected. The controller has a resource budget  $B_h \in \mathbb{N}$  and a probing budget  $B_p \in \mathbb{N}$ . When an arrival is accepted (respectively, probed), we reduce  $B_h$  (respectively,  $B_p$ ) by one.

Formally, we view each time period  $t \in \{T, T-1, ..., 1\}$ as comprising of a mini dynamic program with two stages  $\{t, t-1/2\}$ , driven by external random inputs  $\xi^t \in [n]$  and  $\xi^{t-1/2} \in [n] \times [m]$ . In the first stage t, the controller observes the arriving request  $\xi^t = j$ , and chooses an action in {a,p,r}; in the second stage t-1/2, the reward  $r_{jk}$  (or subtype  $\xi^{t-1/2}=(j,k)\in$  $[n] \times [m]$ ) is drawn with probability  $q_{ik}$ , and the available actions are  $\{a, r\}$  if the first-stage action is p, and  $\emptyset$ otherwise. We augment the state space with a variable that captures the first-stage decision (i.e., whether we accept/reject without probing or probe). The state space S of the controlled process is thus S = $\{(b_h, b_p, \diamond) : b_h, b_p \in \mathbb{N}, \diamond \in \{a, p, r, \emptyset\}\}, \text{ where } b_h, b_p \text{ are }$ the residual hiring and probing budgets. In first stage of each period, we set  $\diamond = \emptyset$ , and only collect rewards in the second stage in each period. See Figure 3 for an illustration.

# **4.1. Offline Benchmark and Online Policy for Probing** We now apply the RABBI framework for online probing.

**4.1.1. Offline Benchmark.** We define Offline to be the controller that knows the public types of all arrivals in advance (i.e., it knows  $Z_j^t$ , the number of type-j items that will arrive in the last t periods), but does not know the realization of the rewards (subtypes). Formally, Offline is endowed with the canonical filtration given by  $\Theta = [T]$  and  $G_\theta = \xi^\theta$  (see Definition 1): with t steps to go, Offline has the information filtration  $\mathcal{G}_t = \sigma(\{\xi^t : t \in [T]\} \cup \{\xi^\tau : \tau \geq t\})$ . Note that since Offline does not know the actual rewards, it still needs to solve a dynamic program to decide whether to probe an arrival.

**Figure 3.** Actions/Transitions in Online Probing in Periods t, t-1/2, and t-1, with Inputs  $\xi^t = j$  and  $\xi^{t-1/2} = R^j$ 



*Notes.* Numbers below the arrows represent the reward of a transition. At t, available actions are  $\{a, p, r\}$  (i.e., accept, probe, reject; from top to bottom). At t - 1/2, if we chose to probe in the first stage (i.e., are in the middle state), then available actions are  $\{a, r\}$ .

**4.1.2. Relaxed Value Function.** Since solving for Off-LINE's optimal actions may be nontrivial, we next construct a relaxed value function  $\varphi$ , using the following LP parametrized by  $(b_h, b_p, \mathbf{z}) \in \mathbb{N}^2 \times \mathbb{R}^n_{>0}$ :

$$P[b_{h}, b_{p}, \mathbf{z}])$$
maximize: 
$$\sum_{j,k} r_{jk} x_{jka} + \sum_{j} \bar{r}_{j} x_{ja},$$
subject to: 
$$\sum_{j,k} x_{jka} + \sum_{j} x_{ja} \leq b_{h},$$

$$\sum_{j} x_{jp} \leq b_{p},$$

$$x_{ja} + x_{jp} + x_{jr} = z_{j} \qquad \forall j \in [n],$$

$$x_{jka} + x_{jkr} = q_{jk} x_{jp} \quad \forall j \in [n], k \in [m],$$

$$\mathbf{x} \geq 0.$$

$$(7)$$

Intuitively,  $P[b_h, b_p, \mathbf{z}]$  can be understood as follows: given current resource and probing budgets  $\mathbf{b}$  and future arrivals  $\mathbf{z}$ , the decision variables  $\mathbf{x} \in \mathbb{R}^{3n+2nm}_{\geq 0}$  represent action summaries, where  $x_{ja}, x_{jx}, x_{jp}$  are the total number of future type-j arrivals that are accepted without probing, rejected without probing, and probed, respectively, and  $x_{jka}, x_{jkr}$  are the number of probed future type-j arrivals that are revealed to have reward  $r_{jk}$ , and then accepted/rejected, respectively. The first two constraints implement the resource budget and probing budget; the third ensures the number of type-j items accepted, probed, or rejected equals arrivals of that type. Finally, the last constraint guarantees that a  $q_{jk}$  fraction of probed type-j items have subtype k (i.e., reward  $r_{jk}$ ).

To construct relaxed value  $\varphi$ , recall that a state is of the form  $s = (b_h, b_p, \diamond)$  with  $\diamond \in \{a, p, r, \emptyset\}$ . For period t (i.e., first stage,  $\diamond = \emptyset$ ), we define  $\varphi(t, (b_h, b_p, \emptyset)|\mathcal{G}_t) := P[b_h, b_p, Z^t]$ . For t - 1/2 (i.e., second-stage decisions), we modify  $\varphi$  to incorporate the action (a, p, r) taken in

the first stage. Overall, our relaxation is defined as follows:

$$\varphi(t - 1/2, (b_h, b_p, \diamond)|\mathcal{G}_t) 
= \begin{cases} r_{\xi^{t-1/2}} + P[(b_h, b_p), Z^{t-1}] & \diamond = \mathbf{a} \\ \max\{r_{\xi^{t-1/2}} + P[(b_h - 1, b_p), \\ Z^{t-1}], P[(b_h, b_p), Z^{t-1}]\} & \diamond = \mathbf{p} \\ P[(b_h, b_p), Z^{t-1}] & \diamond = \mathbf{r}. \end{cases}$$
(8)

### 4.1.3. Value Function Estimate and Online Policy.

Finally, we can use the relaxed value function  $\phi$  in Section 8 to construct an estimated value function  $\hat{\varphi}$ by replacing  $Z^t$  with  $\mathbb{E}_{\xi^{t-1/2}}[Z^t]$ . Using this, we get our online policy specified in Algorithm 2.

### Algorithm 2 (Probing RABBI)

**Input:** Access to solutions of  $(P[\mathbf{b}, \mathbf{z}])$ 

**Output:** Sequence of decisions for Online.

- 1: Initialize budgets  $(B_h^T, B_p^T) \leftarrow (B_h, B_p)$
- 2: **for** period t = T, ..., 1 **do**
- Compute  $X^t$ , an optimal solution to  $(P[B^t, \mathbb{E}[Z^t]])$ .
- Observe the arrival, say it is of type j, then take 4: action  $\hat{U}^t \in \arg\max_{u=a,p,r} \{X_{j,u}^t\}.$
- If  $\hat{U}^t = \mathbf{r}$  or  $\hat{U}^t = \mathbf{a}$ : collect zero or random  $R_i$ , 5: respectively.
- If  $\hat{U}^t = p$ : probe the arrival to observe  $R_i = r_{jk}$ , 6: then take action  $\arg\max_{u=a,r} \{X_{j,k,u}^t\}$ . Update budgets  $B^{t-1}$  accordingly.
- 7:

Remark 2 (Probing Cost). Our approach can also handle a setting where the controller has no probing budget, but instead incurs a penalty  $c_i$  when probing a type-iarrival. The only change to results and proofs is in the definition of  $P[\mathbf{b}, Z]$ , where we drop the constraint involving the probing budget, and modify the objective to be  $\max\{\sum_{i,k} r_{ik} x_{ika} + \sum_i \bar{r}_i x_{ia} - \sum_i c_i x_{ip}\}.$ 

### 4.2. Regret Analysis for Online Probing

We now provide a brief outline of the proof of Theorem 2, which guarantees that Algorithm 2 has a regret that is independent of T,  $B_h$  and  $B_p$ . Complete proofs are provided in online Appendix EC.2.

The main part of the proof involves showing that  $\varphi$ as defined in Equation (8) obeys the Bellman inequalities (Definition 2) with appropriately chosen Bellman loss. The first ingredient for this is provided by the following lemma, which establishes initial ordering for our relaxed value  $\varphi$ .

**Lemma 2.** For any  $b_h, b_p \in \mathbb{N}$ , and arrivals Z,  $\mathbb{E}[V(T,$  $(b_h, b_p)|\mathcal{G}_T)$ ]  $\leq \mathbb{E}[\varphi(T, (b_h, b_p, \emptyset)|\mathcal{G}_T)].$ 

This follows from a standard argument, where we argue that any offline policy induces action summaries that satisfy the constraints defining  $\varphi$ . The proof is provided in online Appendix EC.2.

The bulk of the work is in establishing monotonicity, which we do via the following lemma. Recall the definitions of exclusion sets, satisfying actions and maximum loss (Definitions 3, 4, and 5).

**Lemma 3.** Let  $\bar{X}$  be a maximizer of  $(P[(b_h, b_n), Z^t])$  for some period t, and suppose  $\xi^t = i$ . Then we have the following implications for satisfying actions:

- 1. If  $\bar{X}_{ia} \geq 1$ , then accepting at time t is a satisfy-
- 2. If  $\bar{X}_{ir} \geq 1$ , then rejecting at time t is a satisfying action.
- 3. If  $\bar{X}_{ip} \ge 1$ , and  $\xi^{t-1/2} = (i,k)$  is such that either  $\bar{X}_{ika} \ge 1$  or  $\bar{X}_{ikx} \ge 1$ , then probing at time t, followed by accepting (if  $\bar{X}_{ika} \ge 1$ ) or rejecting (if  $\bar{X}_{ikr} \ge 1$ ) at time t - 1/2 is a satisfying action.

Finally  $\varphi$  satisfies the Bellman inequalities with Bellman  $loss L_B(t, (b_h, b_p)) = r_{\varphi} \mathbb{1}_{\mathcal{B}(t, b_h, b_n)}$ , where  $\mathcal{B}$  are exclusion sets defined as:

$$\mathcal{B}(t,b_h,b_p) = \{ \omega \in \Omega : \not\exists \bar{X} \text{ solution to } (P[(b_h,b_p),Z^t])$$
s.t. (1) or (2) or (3) hold\}.

The proof generalizes the argument in Example 2 for Online Knapsack. We provide a brief outline here, and defer the details to online Appendix EC.2. First, observe that the monotonicity condition in Definition 2 translates to the following condition in the online probing setting:

$$\varphi(t, (b_h, b_p, \emptyset)|\mathcal{G}_t)$$

$$\leq \max_{\diamond \in \{a, p, r\}} \{ \mathbb{E}_{\xi^{t-1/2}} [\varphi(t - 1/2, (s_\diamond, \diamond)|\mathcal{G}_{t-1/2})|\mathcal{G}_t] \}$$

$$\forall \omega \notin \mathcal{B}(t, b_h, b_p),$$

where the state  $s_{\diamond} = (b_h - 1, b_p)$  if  $\diamond = a$ ,  $s_{\diamond} = (b_h, b_p - 1)$ if  $\diamond = p$ , and  $s_{\diamond} = (b_h, b_p)$  if  $\diamond = r$ . Moreover, given  $\xi^t = i$ , we have from Equation (8) that  $\mathbb{E}_{\xi^{t-1/2}}[\varphi(t-1/2,$  $(s_{\diamond},\diamond)|\mathcal{G}_{t-1/2}|\mathcal{G}_t| = P[(b_h,b_p),Z^{t-1}] \text{ if } \diamond = \mathtt{r}, \text{ and } r_{\xi^{t-1/2}} +$  $P[(b_h - 1, b_p), Z^{t-1}]$  if  $\diamond =$  a. Now for cases (1) and (2), the claim in the lemma follows directly by invoking Lemma 1. Finally, case (3) (where  $X_{ip} \ge 1$ ) also follows from using Lemma 1, but in a somewhat more technical way (see online Appendix EC.2 for details).

Using Lemmas 2 and 3, we can complete the regret analysis for Algorithm 2.

**Proof of Theorem 2.** By Theorem 5, we have that Regret  $\leq r_{\varphi} \sum_{t} (\mathbb{1}_{\mathcal{B}(t,S^t)} + \mathbb{1}_{\mathcal{Q}(t,S^t)})$ . To bound this, we proceed in two steps: bounding the measure of the exclusion sets  $\mathcal{B}$ , and the disagreement sets  $\mathcal{Q}$ . We conclude using the fact that  $r_{\varphi} \leq \max_{j,k} r_{jk}$ .

To bound the measure of the exclusion sets  $\mathcal{B}$ , let Xbe the solution to  $(P[\mathbf{b}, Z^t])$ , and note that Lemma 3 guarantees that there is zero Bellman loss if (1)  $\max\{\bar{X}_{ja}, \bar{X}_{jr}\} \ge 1$ , or (2)  $\bar{X}_{jp} \ge 1$  and  $\max\{\bar{X}_{jka}, \bar{X}_{jkr}\} \ge 1$ . The exclusion set  $\mathcal{B}(t, \mathbf{b})$  comprises sample paths where both (1) and (2) fail.

Note that any feasible solution to  $(P[\mathbf{b}, Z^t])$  satisfies  $x_{ja} + x_{jp} + x_{jr} = Z_j^t \ \forall j$  and  $x_{jka} + x_{jkr} = q_{jk}x_{jp} \ \forall j,k$ . If  $Z_j^t \geq 3$ , then one of the variables  $x_{ja}, x_{jp}, x_{jr}$  must be at least 1. On the other hand, we need  $q_{jk}x_{jp} \geq 2$  to guarantee that one of  $x_{jka}, x_{jkr}$  is at least 1. Thus we have

$$\mathbb{P}\left[\mathcal{B}(t,b)|\xi^{t-1/2} = (j,k)\right] \le \mathbb{P}\left[Z_j^t < \frac{6}{q_{jk}}\right]$$

$$= \mathbb{P}\left[Z_j^t - \mu_j(t) < -\mu_j(t)\left(1 - \frac{6}{\mu_j(t)q_{jk}}\right)\right].$$
(9)

Restricting  $\mu_j(t) \ge 12/q_{jk}$  to ensure the RHS of Equation (9) is positive, we can use a standard Chernoff bound (see Boucheron et al. 2013) to get  $\mathbb{P}[\mathcal{B}(t,b)|\xi^{t-1/2}=(j,k)] \le e^{-2(p_j/2)t} + \mathbb{1}_{\{t \le 12/(p_iq_{jk})\}}$ . Finally,

$$\sum_{t} \mathbb{P}[\mathcal{B}(t, B^{t})] \leq \sum_{t} \sum_{j} p_{j} e^{-2(p_{j}/2)t}$$

$$+ \sum_{t} \sum_{j,k} p_{j} q_{jk} \mathbb{1}_{\{t \leq 12/(p_{j}q_{jk})\}}$$

$$\leq \sum_{j} \frac{2}{p_{j}} + 12.$$

To bound the information loss  $\Sigma_t \mathbb{P}[\mathcal{Q}(t,S^t)]$ , recall  $\mathcal{Q}(t,S^t) \subseteq \Omega$  is the event where  $\hat{U}^t$  is not satisfying. Let  $\bar{X}$  be a solution to  $(P[\mathbf{b},Z^t])$ , t a first stage, and let  $j=\xi^t$ . We now have two cases depending on if  $\hat{U}^t \in \{\mathbf{a},\mathbf{r}\}$  or  $\hat{U}^t=\mathbf{p}$ . First, if  $\hat{U}^t \in \{\mathbf{a},\mathbf{r}\}$ , then according to Lemma 3, accepting or rejecting is satisfying whenever  $\max\{\bar{X}_{j\mathbf{a}},\bar{X}_{j\mathbf{r}}\} \geq 1$ . Since  $X^t(\xi^t,\hat{U}^t)=\max\{X^t(\xi^t,u):\}$   $u=\mathbf{a},\mathbf{p},\mathbf{r}$  and  $X^t_{j\mathbf{a}}+X^t_{j\mathbf{p}}+X^t_{j\mathbf{r}}=\mu_j(t)$ , we have

$$\mathbb{P}\left[\bar{X}(j,\hat{U}^t) < 1|X^t(j,\hat{U}^t) \ge \mu_j(t)/3\right]$$
  
$$\le \mathbb{P}\left[\left\|\bar{X} - X^t\right\|_{\infty} \ge \mu_j(t)/3\right].$$

On the other hand, if  $\hat{U}^t = p$ , the error is bounded by

$$\begin{split} \mathbb{P}\left[\bar{X}_{jp} < 1 \text{ or } \bar{X}_{\xi^{t-1/2},u} < 1 \middle| X_{jp}^{t} \geq \frac{\mu_{j}(t)}{3}, X_{\xi^{t-1/2},u}^{t} \\ \geq \frac{q_{\xi^{t-1/2}}\mu_{j}(t)}{6} \right] \leq \mathbb{P}\left[ \left\| \bar{X} - X^{t} \right\|_{\infty} \geq \frac{q_{\xi^{t-1/2}}\mu_{j}(t)}{6} \right], \end{split}$$

where u is the action with largest value between the variables  $X^t(\xi^{t-1/2}, \mathbf{a}), X^t(\xi^{t-1/2}, \mathbf{r})$ .

Thus, regardless of the action  $\hat{U}^t$ , the probability of choosing a nonsatisfying action is bounded by  $\mathbb{P}[\|\bar{X} - X^t\|_{\infty} \ge \min_k q_{jk} \cdot \mu_j(t)/6]$ . Moreover, standard LP sensitivity results (Mangasarian and Shiau 1987, theorem 2.4) imply that there exists  $\kappa$  depending on  $\mathbf{q}$ , n, m alone, such that  $\|\bar{X} - X^t\|_{\infty} \le \kappa \|Z^t - \mu(t)\|_1$ . Finally, the measure of sets  $\mathcal{Q}$  where Online chooses a nonsatisfying action is bounded by

$$\sum_{t} \mathbb{P}[\mathcal{Q}(t, S^{t})]$$

$$\leq \sum_{t} \mathbb{P}\left[\left\|Z^{t} - \mu(t)\right\|_{1} \geq \min_{k} q_{jk} \cdot \mu_{j}(t) / 6\kappa\right] < \infty.$$

The summability follows arguments presented in (Vera and Banerjee 2020), based on standard concentration bounds.

### 5. Dynamic Pricing

We now apply our framework to dynamic pricing. In the basic setting, we have d resources and n customer types. Each customer type has a private reward for a set of resources. The controller observes the customer type, and if the corresponding set of resources is available, posts a price. The customer then purchases iff the requested set is available and the posted price is below the private reward. The resource consumption is encoded in a matrix  $A \in \{0,1\}^{d \times n}$ . In Section 5.5, we generalize to settings where rather than requesting a specific set of products, customers choose between multiple substitute bundles of resources.

We consider the following formal model: at time t, type  $j \in [n]$  arrives with probability  $p_j$ , is seen by the controller, who then posts a price  $f_{jl}$  from a set of available prices  $\{f_{j1}, \ldots, f_{jm}\}$ . The customer then draws a private reward  $R^t \sim F_j$ , and a purchase occurs iff  $R^t > f_{jl}$ . If the customer buys,  $f_{jl}$  is collected and the inventory decreases by  $A_j$ . On the other hand, if the customer does not buy, the controller collects zero and the inventory remains unchanged.

### 5.1. Offline Benchmark and Online Policy for Dynamic Pricing

5.1.1. Offline Benchmark. Note that for each customer type j, there are  $Z_i^T$  arrivals, and hence  $Z_i^T$  draws from the distribution  $F_i$ . We now define our benchmark by considering Offline to be a controller that knows the realized histogram of these draws, that is, for each j, Offline knows the empirical distribution of the  $Z_i^T$  rewards. Moreover, at the end of each period t, Offline also observes the realized valuation  $R^t$  whether there is a sale. Note that Offline does not know the exact sequence of these rewards, and so is not a fullinformation benchmark. For example, say  $Z_1^T = 15$ and we reveal that 10 arrivals type-1 have private reward \$1 and five arrivals type-1 have private reward \$2. Now, upon observing a type-1 arrival, Offline concludes that the reward is \$2 with probability  $\frac{5}{15}$ . Now if the arrival had value \$1, then, the next time Offline observes a type-1, its belief is that the reward is \$2 with probability  $\frac{5}{14}$ .

Formally, for each j, suppose the prices are ordered  $f_{j1} > f_{j2} > \ldots > f_{jm}$ . Denote  $\xi^t \in [n]$  to be the type of the arrival at time t. To define Offline, we introduce a sequence of independent random vectors  $\{Y^t: t=T, T-1,\ldots,1\}$  where  $Y^t_{jl}:=\mathbb{1}_{\{\xi^t=j,R^t>f_{jl}\}}$ ; in other words,  $Y^t_{jl}$  is the indicator of whether a price  $f_{jl}$  or lower is accepted by the type-j at time t. We define  $Q_{jl}(t):=\frac{1}{Z^t_i}\sum_{\tau=1}^t Y^\tau_{jl}$  to be the fraction of type-j

customers who accept price  $f_{jl}$  in the last t periods. Observe that  $Q_{jl}(t)$  is a martingale with  $\mathbb{E}[Q_{jl}(t)] = \bar{F}_j(f_{jl})$  and  $Q_{jl}(t) = \frac{Z_j^{t+1}}{Z_j^t}Q_{jl}(t+1) - \frac{1}{Z_i^t}Y_{jl}^{t+1}$ .

Offline's information is now given by the filtration  $\mathcal{G}_t = \sigma(\{Q(\tau), Z^\tau : \tau \geq t\})$ , that is, at every time t, Offline knows the total demand  $Z_j^t$  and the empirical averages  $Q_{jl}(t)$ , but not the sequence of rewards. This coincides with the canonical filtration (Definition 1) with variables  $(Q_{jl}(T), Z_j^T : j \in [n], l \in [m])$ . The filtration  $\mathcal G$  is strictly coarser than the full-information filtration, which would correspond to revealing all the variables  $Y^T, Y^{T-1}, \ldots, Y^1$  instead of their empirical averages.

# **5.1.2. Relaxed Value Function.** Consider the following LP, parameterized by (**b**, **q**, **z**):

$$(P[\mathbf{b}, \mathbf{q}, \mathbf{z}]) \quad \text{maximize:} \quad \sum_{j,l} f_{jl} q_{jl} x_{jl}$$

$$\text{subject to:} \quad \sum_{j,l} a_{ij} q_{jl} x_{jl} \le b_i \quad \forall i \in [d]$$

$$\sum_{j,l} x_{jl} + x_{jr} = z_j \quad \forall j \in [n]$$

$$\mathbf{x} \ge 0$$

$$(10)$$

We define the relaxed value as  $\varphi(t, \mathbf{b}|\mathcal{G}_t) := P[\mathbf{b}, Q(t), Z^t]$ , and the corresponding estimated value as  $\hat{\varphi}(t, \mathbf{b}) := P[\mathbf{b}, \mathbf{q}, t\mathbf{p}]$ , where  $q_{jl} = \bar{F}_j(f_{jl})$ . The resulting RABBI policy is presented in Algorithm 3.

### Algorithm 3 (Pricing RABBI)

**Input:** Access to solutions of  $(P[\mathbf{b}, \mathbf{q}, \mathbf{z}])$ 

**Output:** Sequence of decisions for Online.

- 1: Set  $B^T \leftarrow B$  as the given initial budget and  $q_{il} \leftarrow \bar{F}_i(f_{il})$
- 2: **for** t = T, ..., 1 **do**
- 3: If the arrival is type j and  $A_j \nleq B^t$ : not enough resources, reject and go to t-1.
- 4: Compute  $X^t$ , an optimal solution to  $(P[B^t, \mathbf{q}, t\mathbf{p}])$ .
- 5: Let  $l \in \arg \max\{X_{j,l}^t : l = 1, ..., m, r\}$ . If l = r, reject and go to t 1. Else post price  $f_{jl}$ .
- 6: If  $R^t > f_{jl}$ , collect  $f_{jl}$  and  $B^{t-1} \leftarrow B^t A_j$ ; else  $B^{t-1} \leftarrow B^t$ .

To get some intuition into the LP ( $P[\mathbf{b}, \mathbf{q}, \mathbf{z}]$ ), note that if  $q_{jl} = \bar{F}_j(f_{jl})$ , that is, the probability that price  $f_{jl}$  is accepted by a type-j customer and  $z_j$  is the number of type-j arrivals, then ( $P[\mathbf{b}, \mathbf{q}, \mathbf{z}]$ ) can be interpreted as follows: the variable  $x_{jl}$  represents the number of times that price  $f_{jl}$  is offered, with  $\sum_{j,l} f_{jl}q_{jl}x_{jl}$  the expected reward from the corresponding arrivals. Each time price  $f_{jl}$  is offered,  $a_{ij}q_{jl}$  units of resource i are consumed in expectation, and hence  $\sum_{j,l} a_{ij}q_{jl}x_{jl}$  is the total expected consumption of resource i. Finally, at

most one price is offered per arrival, which is captured by  $\sum_{l} x_{jl} + x_{jr} = z_{j}$ , where  $x_{jr}$  is the number of rejected type-j customers.

### 5.2. Bellman Inequalities and Bellman Loss

We first argue that our choice of  $\phi$  satisfies the Bellman inequalities.

**Lemma 4.** Let  $V(T, B|\mathcal{G}_T)$  be the value of Offline's optimal policy, and  $\varphi(t, \mathbf{b}|\mathcal{G}_t) = P[\mathbf{b}, Q(t), Z^t]$  be the relaxed value with optimal solution X:

- 1. The relaxed value satisfies the initial ordering condition:  $\mathbb{E}[V(T, B|\mathcal{G}_T)] \leq \mathbb{E}[\varphi(T, B|\mathcal{G}_T)]$ .
- 2. If the arriving type is j and  $\max_{l}\{X_{jl}\} \ge 1$ , then  $\mathbb{E}[L_B(t, \mathbf{b})] \le 0$ .
- 3. If the arriving type is j and  $X_{jl} \ge 1$ , then posting  $f_{jl}$  is a satisfying action.

We omit the proof of the initial ordering in item (1), as it is similar to that of Lemma 2. Next we present the main ingredients for obtaining the monotonicity property (items (2) and (3)). Complete details are deferred to online Appendix EC.3. For ease of exposition, when the controller rejects, the controller can equivalently post  $f_{jr} = \infty$  such that  $\bar{F}_j(f_{jr}) = 0$  with the convention  $0 \times \infty = 0$ .

We start by recalling the monotonicity condition (Definition 2). Denote  $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot|\mathcal{G}_t]$ . If the inventory is  $\mathbf{b} \geq A_j$ , the random reward of posting price  $f_{jl}$  at t is  $f_{jl}Y_{jl}^t$  and the random new inventory is  $\mathbf{b} - A_jY_{jl}^t$ , thus monotonicity corresponds to:

$$\varphi(t+1,\mathbf{b}) \leq \max_{l \in [m] \cup \{\mathbf{r}\}} \left\{ \mathbb{E}_{t+1} \left[ f_{jl} Y_{jl}^{t+1} + \varphi \left( t, \mathbf{b} - A_j Y_{jl}^{t+1} \right) \right] \right\}$$
$$+ \mathbb{E}_{t+1} \left[ L_B(t+1,\mathbf{b}) \right].$$

Because Q is a martingale, we have  $\mathbb{E}_t[Y^t] = Q(t)$ , and we can further simplify the condition to:

$$\varphi(t+1,\mathbf{b}) \leq \max_{l \in [m] \cup \{r\}} \left\{ f_{jl} Q_{jl}(t+1) + \mathbb{E}_{t+1} \left[ \varphi(t,\mathbf{b} - A_j Y_{jl}^{t+1}) \right] \right\} \\
+ \mathbb{E}_{t+1} [L_B(t+1,b)]. \tag{11}$$

Define  $L_B(t+1,b,j,l) := \varphi(t+1,\mathbf{b}) - f_{jl}Q_{jl}(t+1) - \mathbb{E}_{t+1}$   $[\varphi(t,\mathbf{b}-A_jY_{jl}^{t+1})]$ , which corresponds to the loss in Equation (11) when we assume a specific price  $f_{jl}$  is posted. Recall we define  $\varphi(t+1,\mathbf{b}) = P[\mathbf{b},Q(t+1),Z^{t+1}]$ . Moreover, for an arrival of type j and any solution X of  $P[\mathbf{b},Q(t+1),Z^{t+1}]$ , if  $X_{jl} \geq 1$ , then using Lemma 1, we have  $P[\mathbf{b},Q(t+1),Z^{t+1}] = f_{jl}Q_{jl}(t+1) + P[\mathbf{b}-A_jQ_{jl}(t+1),Q(t+1),Z^t]$ . Thus, assuming  $X_{jl} \geq 1$ , we can write the loss in the Bellman inequality as

$$L_{B}(t+1,\mathbf{b},j,l) = P[\mathbf{b} - A_{j}Q_{jl}(t+1), Q(t+1), Z^{t}]$$
$$-\mathbb{E}_{t+1}[P[\mathbf{b} - A_{j}Y_{jl}^{t+1}, Q(t), Z^{t}]]. \quad (12)$$

Observe that  $L_B(t, \mathbf{b}, j, l)$  is characterized by a random LP that depends on  $Y^{t+1}$  (which is unknown at time t+1), see Equation (12). To complete item (2) of Lemma 4, it remains to prove that  $L_B(t, \mathbf{b}, j, l)$  characterized in (12) satisfies  $\mathbb{E}_t[L_B(t, \mathbf{b}, j, l)] \leq 0$ . This is proved in online Appendix EC.3 by arguing that the term in (12) is upper bounded by a zero-mean random variable.

We can then conclude that, for each l with  $X_{jl} \ge 1$ ,  $\mathbb{E}_{t+1}[L_B(t+1,b,j,l)] \le 0$ , so that  $\varphi(t+1,\mathbf{b}) \le \mathbb{E}_{t+1}[f_{jl}Q_{jl}(t+1) + \varphi(t,\mathbf{b} - A_jY_{jl}^{t+1})]$ , implying that posting price  $f_{il}$  is a satisfying action, which is item (3) of Lemma 4.

# 5.3. Information Loss and Overall Performance Guarantee

Next we study the disagreement sets  $Q(t, B^t)$  and bound the information loss  $\Sigma_t \mathbb{P}[Q(t, B^t)]$ .

**Proposition 2.** Let X be a solution of  $(P[\mathbf{b}, Q(t), Z^t])$ . If  $X_{jl} \geq 1$ , then posting  $f_{jl}$  is a satisfying action. Furthermore, the information loss is bounded by  $\mathbb{P}[Q(t, B^t)] \leq 1/t^2$  for all  $t \geq c$ , where c depends only on  $(\mathbf{f}, \mathbf{p}, A, F_1, \dots, F_n)$ .

We now give an outline of this proof (for details, refer to online Appendix EC.3). Recall that RABBI chooses l as the maximum entry of the solution to  $(P[\mathbf{b}, \mathbb{E}[Q(t)], \mathbb{E}[Z^t])$ , which is a perturbed version of the object of interest, thus Online needs to guess l such that  $X_{jl} \geq 1$  without the knowledge of Q(t) and  $Z^t$ , creating an information loss.

To build intuition, consider the case where d = 1 and n = 1, that is, selling multiple copies of an item to homogeneous customers. Since there is only one type, we drop the index j. Recall  $f_1 > \ldots > f_m$  and  $q_1 < \ldots < q_m$ . It is easy to check that the solution of  $P[b, \mathbf{q}, t]$  is as follows: (i) if  $b \le tq_1$ , then  $x = (b/q_1, 0, \ldots, 0)$ ; (ii) if  $b > tq_m$ , then  $x = (0, \ldots, 0, t)$ ; (iii) otherwise,

if  $b \in (tq_l, tq_{l+1}]$ , then  $x_{l'} = 0$  for  $l' \neq l, l+1$ , and  $x_l = (tq_{l+1} - b)(q_{l+1} - q_l), x_{l+1} = (b - tq_l)(q_{l+1} - q_l)$ . Figure 4 illustrates this solution, and also shows that for RABBI's guess to be incorrect, Q(t) and  $\mathbb{E}[Q(t)]$  must deviate considerably, which the next lemma indicates is unlikely. This intuition carries over to higher dimensions.

**Lemma 5.** For any  $j \in [n]$ , there is a constant  $c_j$  depending on  $p_j$  only such that, for any time t,  $\mathbb{P}[\max_l Q_{jl}(t) - \mathbb{E}[Q_{jl}(t)] > \sqrt{\frac{\log(t)}{t}}] \leq \frac{c_j}{t^2}$ .

**Proof.** From the DKW inequality (Massart 1990) for empirical measures, we have

$$\mathbb{P}\left[\sup_{l}Q_{jl}(t)-\bar{F}(f_{jl})>\lambda\left|Z^{t}\right|\right]\leq 2e^{-2\lambda^{2}Z_{j}^{t}}.$$

Also for  $Z_j^t \sim \text{Bin}(t, p_j)$ ,  $\mathbb{E}[e^{-\theta Z_j^t}] = (1 - p + pe^{-\theta})^t$ . Setting  $\lambda = \sqrt{\log(t)/t}$ , we get

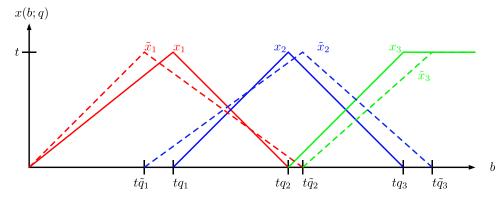
$$\mathbb{P}\left[\sup_{t} Q_{jl}(t) - \bar{F}(f_{jl}) > \sqrt{\frac{\log(t)}{t}}\right]$$

$$\leq 2(1 - p_j + p_j e^{-\theta})^t \quad \text{where } \theta = 2\log(t)/t.$$

Using the inequality  $e^{-\theta} \le 1 - \theta + \theta^2/2$ , an algebraic check confirms the desired inequality.  $\Box$ 

**5.3.1. Stability of Left-Hand Side Perturbations.** As stated in Algorithm 3, Online takes actions based on  $P[\mathbf{b}, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]]$ , whereas Offline uses  $P[\mathbf{b}, Q(t), Z^t]$ . Therefore, for fixed  $(t, \mathbf{b})$ , we need to compare solutions of  $P[\mathbf{b}, \mathbf{q}, \mathbf{z}]$  to those of  $P[\mathbf{b}, \mathbf{q} + \Delta \mathbf{q}, \mathbf{z} + \Delta \mathbf{z}]$ , where  $\Delta$  is the perturbation. Define  $\mathbf{q} = \mathbb{E}[Q(t)]$ ,  $\mathbf{z} = \mathbb{E}[Z^t]$ ,  $\Delta \mathbf{q} = Q(t) - \mathbb{E}[Q(t)]$ , and  $\Delta \mathbf{z} = Z^t - \mathbb{E}[Z^t]$ .

**Figure 4.** (Color online) Solution to the Pricing LP in Equation (10) for the Case d = 1 and n = 1, Which Correspond to Selling Multiple Copies of an Item to Homogeneous Customers



Notes. If  $b/t \in (q_l, q_{l+1}]$ , the prices used by the LP are  $f_l, f_{l+1}$  and the amount of time we offer each is piecewise linear in the budget. For a perturbation  $\tilde{\mathbf{q}}$  of  $\mathbf{q}$ , we superpose the solutions with the different parameters. Our guess is incorrect only when  $\tilde{x}_l \gg 1$  and  $x_l < 1$ , which necessitates a substantial perturbation of  $\mathbf{q}$ .

**Lemma 6** (Selection Program). Let  $V_t = P[\mathbf{b}, \mathbf{q} + \Delta \mathbf{q}, \mathbf{z} + \Delta \mathbf{z}]$  and fix a component (j', l'). Then posting price  $f_{j'l'}$  is satisfying if  $P_S[V_t, \mathbf{q} + \Delta \mathbf{q}, \mathbf{z} + \Delta \mathbf{z}] \geq 1$ , where

$$P_{S}[V_{t}, \mathbf{q} + \Delta \mathbf{q}, \mathbf{z} + \Delta \mathbf{z}]$$

$$:= \max \left\{ x_{j'l'} : \sum_{j,l} f_{jl} (q_{jl} + \Delta q_{jl}) x_{jl} \ge V_{t}, \mathbf{x} \right.$$

$$feasible for P[\mathbf{b}, \mathbf{q} + \Delta \mathbf{q}, \mathbf{z} + \Delta \mathbf{z}] \right\}.$$

In other words,  $Q(t,b,l) = \{\omega \in \Omega : P_S[V_t[\omega], Q(t), Z^t] < 1\}.$ 

**Proof.** This problem selects, among all the solutions of  $P[\mathbf{b}, \mathbf{q} + \Delta \mathbf{q}, \mathbf{z} + \Delta \mathbf{z}]$ , one with the largest component  $X_{j'l'}$ . From Lemma 4, we know that if  $X_{j'l'} \geq 1$ , then posting  $f_{j'l'}$  is satisfying.  $\square$ 

We have converted the condition " $\exists X$  solving  $P[v,\mathbf{q}+\Delta\mathbf{q},\mathbf{z}+\Delta\mathbf{z}]$  with  $X_{j'l'}\geq 1$ " to an optimization program. Let  $\bar{\mathbf{x}}$  be the solution to the proxy  $P[\mathbf{b},\mathbf{q},\mathbf{z}]$  and let  $v_t$  be the objective value (recall that  $V_t$  is the value of  $P[\mathbf{b},\mathbf{q}+\Delta\mathbf{q},\mathbf{z}+\Delta\mathbf{z}]$ ). Since the algorithm picks the price with the largest component, assume  $\bar{x}_{j'l'}=\max_l\bar{x}_{j'l}\gg 1$ . In particular,  $P_S[v_t,\mathbf{q},\mathbf{z}]\gg 1$  for this fixed (j',l'). We want to show that  $P_S[V_t,\mathbf{q}+\Delta\mathbf{q},\mathbf{z}+\Delta\mathbf{z}]\geq 1$  for that particular (j',l'). To that end, we need to bound the difference between  $P_S[V_t,\mathbf{q}+\Delta\mathbf{q},\mathbf{z}+\Delta\mathbf{z}]$  and  $P_S[v_t,\mathbf{q},\mathbf{z}]$ . This difference depends on (i)  $v_t-V_t$ , (ii)  $\Delta$ , and (iii) the dual variables of  $(P_S[V_t,\mathbf{q}+\Delta\mathbf{q},\mathbf{z}+\Delta\mathbf{z}])$ . Observe that the quantities (i)–(iii) are random. We state the result next. The proof is provided in online Appendix EC.3.

**Lemma 7.** There is a constant c that depends only on  $(\mathbf{f}, \mathbf{p}, A, F_1, \dots, F_n)$  such that, for all  $t \ge c$ , with probability  $1 - c/t^2$ ,  $P_S[V_t, Q(t), Z^t] - P_S[v_t, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]] \ge -c\sqrt{t \log(t)}$ .

Lemma 7 leads to the bound in Proposition 2. Indeed, since the LP in Equation (10) has the constraint  $\sum_{l \in [m] \cup \{r\}} \bar{x}_{jl} = tp_j$ , the maximum entry is guaranteed to have a value of at least  $tp_j/(m+1)$ . Therefore, by definition of the selection program,  $P_S[v_t, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]] \ge tp_j/(m+1)$ . We know that posting  $f_{jl'}$  is satisfying whenever  $P_S[V_t, Q(t), Z^t] \ge 1$  (see Lemma 6), hence posting the maximum entry is satisfying provided that  $tp_j/(m+1) - c\sqrt{t\log(t)} \ge 1$ , which holds for all t large enough.

#### 5.4. Numerical Simulations

We test our algorithm on two systems, henceforth the small system and the large system. For each system, we consider a sequence of instances with increasing horizons and initial inventories.

The small system corresponds to the one-dimensional problem (n = 1 and d = 1). In this case, we can solve the DP for small enough horizons and directly compute the optimality gap. The large system corresponds to a

multidimensional problem with n = 20, d = 25, and m = 3. The DP solution is intractable for the large system, yet we can compute the offline benchmark and compare our algorithm against it. The optimality gap, recall, is bounded by the offline versus RABBI gap.

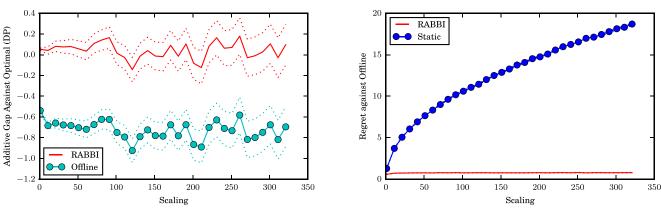
For the small system, the kth instance has budget B = 6k and horizon T = 20k. For each scaling k, we run 100,000 simulations. We consider the following primitives: prices are (1,2,3) and the private reward  $R^t$  has an atomic distribution on (1,2,3) with probabilities (0.3,0.4,0.3). The instance is chosen such that it is dual degenerate for (10), which is supposedly the more difficult case (Jasin 2014). For the large system, the parameters were generated randomly and are reported in online Appendix EC.6. The kth instance has horizon T = 100k and budgets  $B_i = 10k$  for all  $i \in [25]$ .

For the small system, we consider k small enough (short horizon) so that we can compute the optimal policy. This computation becomes intractable already for moderate values of k (rabbi however scales gracefully with k as it only requires resolving an LP in each period). In Figure 5 (left) we display the gap between the optimal solution and both the rabbi and Offline's value. We make two observations: (i) the Offline benchmark outperforms the optimal (as it should), but by a rather small margin; and (ii) rabbi has a constant regret (i.e., independent of k) relative to Offline, and hence constant optimality gap. In contrast, a full-information benchmark would outperform the optimal by too much to be useful.

In Figure 5 (right), we compare RABBI to the optimal static pricing policy, which has regret  $\Omega(\sqrt{k})$  (Gallego and Van Ryzin 1997). In particular, if D(f) denotes the demand at fare f, we choose the static price to be the one that maximizes the revenue function  $f \cdot D(f) = f \cdot T \cdot F(f)$ , subject to the constraint  $D(f) \leq B$ . The solution is the better of two prices: (i) the market clearing price, that is, that satisfies D(f) = B, or (ii) the monopoly price, which maximizes fD(f). We note though that when a continuum of prices is allowed, Jasin (2014) proposes an algorithm (which, like RABBI, is based on resolving an optimization problem in each period) that achieves a regret that is logarithmic in k under certain nondegeneracy assumptions on the optimization problem and differentiability assumptions on the valuation distribution. In contrast, our constant-regret guarantees hold under a finite price menu.

In Figure 6, we display the results for the large system. Here, since the DP is intractable, we use the offline benchmark. The resulting regret is negligible relative to the total value, as captured by the approximation factor on the right-hand side of the figure. We also present the competitive ratio of Offline against the full-information benchmark (this upper bounds the

**Figure 5.** (Color online) Regret in the Small System (n = 1 and d = 1), with Horizon T = 20k and Initial Budget B = 6k, Under Scaling k = 1, 10, 20, ..., 340



Notes. Dotted lines represent 90% confidence interval. (Left) additive gaps against the optimal policy, i.e.,  $V^{DP} - V^{RABBI}$  and  $V^{DP} - V^{OFFLINE}$ . (Right) Regret of two policies against Offline.

competitive ratio of any nonanticipatory policy) and observe that is bounded away from 1, hence showing that the full-information benchmark is  $\Omega(T)$  away from the DP in our randomly generated instance, which confirms the need for our refined benchmark.

### 5.5. Posted Pricing with Customer Choice

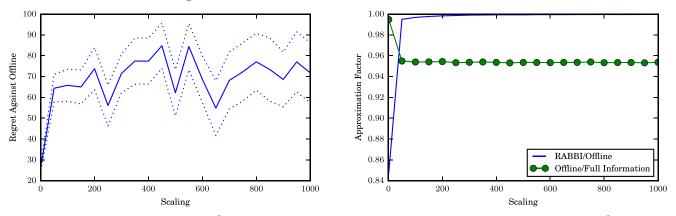
We now consider settings where customers, rather than requesting a specific product, choose between multiple substitutes. As a concrete example, consider a hardware store selling washers and dryers. The store can set a separate price for a washer, a dryer, and also for buying a washer-and-dryer bundle (i.e., one of each). An incoming customer sees the prices and chooses to buy each of the three options (or nothing at all) with some probability depending on the price menu. See Talluri and Van Ryzin (2006, chapter 7) for details on such customer-choice models. For exposition,

we focus here on a single-customer type, with arbitrary (but known) customer-choice model.

As before, the controller chooses a price to post for each product and selling one unit of product  $j \in [n]$  depletes resources according to  $A_j \in \{0,1\}^d$ . There is a discrete set of assortment menus, denoted by  $\mathcal{A}$ . An assortment  $\alpha \in \mathcal{A}$  is associated with a vector of prices  $(f_{1\alpha}, \ldots, f_{n\alpha})$ , one price per product. Setting  $f_{j\alpha} = \infty$  corresponds to not offering product j. Note that if each product's price is restricted to take one of m distinct values, then there are at most  $\mathcal{A} \leq m^n$  different assortments. The actual number of relevant assortments might, however, be much smaller than this.

An arriving customer, when offered assortment  $\alpha$ , chooses to buy product j with a probability  $p_j(\alpha)$ , with  $\sum_{j=0}^{n} p_j(\alpha) = 1$  (where we use j = 0 for the nopurchase option). These probabilities might be derived, for example, from a standard family such as

**Figure 6.** (Color online) Performance in the Large System (n = 20 and d = 25) with Horizon T = 100k and Initial Budgets  $B_i = 10k$  for  $i \in [25]$ , Under Scaling k = 1, 2, ..., 1,000



Notes. (Left) Regret against Offline, that is,  $V^{\text{Offline}} - V^{\text{Rabbl}}$ ; dotted lines represent 90% CI. (Right) Approximation factors  $V^{\text{Rabbl}}/V^{\text{Offline}}$  of Rabbi compared with Offline, and  $V^{\text{Offline}}/V^{\text{Full-Info}}$  of Offline against the full-information benchmark. This shows that the full-information benchmark is indeed too loose, as it is  $\Omega(T)$  away from the DP.

the multinomial-logit model, nested logit model, and so on; our results do not need any specific structure on the choice probabilities (although assuming more structure may lead to better regret scaling with respect to the number of price menus and more efficient ways of solving the resulting LP relaxation).

The process unfolds as follows: (i) at time t, the controller posts an assortment  $\alpha \in \mathcal{A}$ ; (ii) with probability  $p_j(\alpha)$ , the arriving customer buys one unit of product j (with product 0 corresponding to no purchase). Now given the choice probabilities, we can simulate the choice model as follows: we assume w.l.o.g. that the customer arriving at time t is endowed with an i.i.d. random variable  $\xi^t \sim \text{Uniform}(0,1)$ , and assert that the customer buys product j if  $\xi^t \in [\sum_{j'=0}^{j-1} p_{j'}(\alpha), \sum_{j'=0}^{j} p_{j'}(\alpha)]$ . Note that the order of products here is arbitrary.

Applying RABBI to this setting gives the following result.

**Theorem 6** (Dynamic Pricing with Customer Choice). For any choice model with probabilities and prices  $(p_j(\alpha), f_{j\alpha}: j \in [n], \alpha \in \mathcal{A})$ , RABBI obtains a regret that depends only on (A, p, f), but is independent of the horizon length T and initial budget levels  $B \in \mathbb{N}^d$ .

**5.5.1. Algorithm and Analysis.** The following LP extends Equation (10) to incorporate consumer choice:

$$\begin{split} \left(P\big[\mathbf{b},\mathbf{q},\mathbf{z}\big]\right) \quad \text{maximize:} \quad & \sum_{\alpha \in \mathcal{A}} x_{\alpha} \sum_{j \in [n]} f_{j\alpha} q_{j\alpha} \\ \text{subject to:} \quad & \sum_{\alpha \in \mathcal{A}} \sum_{j \in [n]} a_{ij} q_{j\alpha} x_{\alpha} \leq b_{i} \quad \forall \, i \in [d], \\ & \sum_{\alpha \in \mathcal{A}} x_{\alpha} = t \\ & \mathbf{x} \geq 0. \end{split}$$

Here,  $q_{j\alpha}$  stands for the fraction of customers that would buy product j if presented with the price assortment  $\alpha$ . RABBI re-solves, in each period, this LP with the expected fraction  $q_{j\alpha} = p_j(\alpha)$ . In contrast, Offline knows  $Q_{j\alpha}(t)$ , the realized fraction of customers that, given assortment  $\alpha$ , would buy product j (formally, Offline is equipped with the canonical augmented filtration with variables  $(Q_{j\alpha}(T):j\in [n],\alpha\in\mathcal{A})$ ), and solves Equation (13) with  $q_{j\alpha}=Q_{j\alpha}(t)$ , where:

$$\begin{split} Q_{j\alpha}(t) := & \frac{1}{t} \sum_{\tau=1}^{t} Y_{j\alpha}^{t} \\ \text{where } Y_{j\alpha}^{t} := & \mathbb{1}_{\{\sum_{j'=0}^{j-1} p_{j'}(\alpha) \leq \xi^{t} \leq \sum_{j'=0}^{j} p_{j'}(\alpha)\}}. \end{split}$$

With the (re)defined key ingredients—namely the LP in Equation (13) and Offline's information structure—it is evident that that the analysis of this expanded

model is identical to that of the basic (no-choice) pricing setting with obvious changes. For example, if assortment  $\alpha$  is posted at time t, the random collected reward is  $\sum_j Y^t_{j\alpha} f_{j\alpha}$  and the random inventory at t-1 is  $b-\sum_j A_j Y^t_{j\alpha}$ . In turn, the Bellman loss in Equation (12) takes on the following form:

$$L_{B}(t+1,\mathbf{b},\alpha) = P\left[\mathbf{b} - \sum_{j} A_{j}Q_{j\alpha}(t+1), Q(t+1), t\right]$$
$$-\mathbb{E}_{t+1}\left[P\left[\mathbf{b} - \sum_{j} A_{j}Y_{j\alpha}^{t+1}, Q(t), t\right]\right].$$
(14)

Now we have a sum over products *j*, but the analysis goes through via linearity of expectations.

**5.5.2. Numerical Simulations.** We demonstrate our algorithm for the following simple choice model with two resources (R1, R2), and three products (R1), R2), R1, R2) (for example, a hardware store selling washers (R1), dryers (R2), or washer-and-dryer combos (R1, R2). The controller has initial inventories of each resource, and can choose among one of seven price assortments: high and low prices with/without discounts for buying the bundle, and price menus assuming stock-out of either or both resource. The price menus and choice probabilities are detailed in Table 1. We run RABBI for this instance while scaling the horizon and initial inventory (see Figure 7).

### 6. Online Knapsack with Distribution Learning

Finally, we consider the distribution-agnostic online knapsack setting. We study first the full feedback setting, and in Section 6.3 extend to censored feedback. As in the baseline OnlineKnapsack, at each time t, the arrival is of type  $j \in [n]$  with known probability  $p_j$ . Type j has a known weight  $w_j$  and random reward  $R_j$ , drawn from a distribution  $F_j$ , with  $r_j := \mathbb{E}[R_j]$ . Critically, we assume  $r_j$  and  $F_j$  are unknown to Online.

The reward  $R_j$  is revealed only after the decision to accept/reject has been made. At the end of each period, we observe the realization of both accepted and rejected items. In contrast, Offline has access to the distribution  $F_j$ , but not to the realizations. We assume that, before the process starts, we are given one sample of each type, and with t periods to go, define  $R_j^t$  to be the empirical average of the observed rewards for type-j arrivals.

As in probing, we divide each period  $t \in \{T, T-1,\ldots,1\}$  into two stages, t and t-1/2. In the first stage (i.e., period t), the input reveals the type  $j \in [n]$ , and in second stages (i.e., period t-1/2), the reward is revealed. The random inputs are given by  $\xi^t \in [n]$  and  $\xi^{t-1/2} \in \mathbb{R}$ . The state space is  $\mathcal{S} = \mathbb{R}_{\geq 0} \times \{\emptyset, \mathsf{a}, \mathsf{r}\}$ ,

**Table 1.** Example with Seven Assortments

Parameter	Products	High	High-discount	Low	Low-discount	Only R2	Only R1	Stock-out
$f_{j\alpha}$ $p_i(\alpha)$	{R1} {R2} {R1,R2} {R1}	5 5 10 0.2	5 5 9 0.2	3 3 6 0.3	3 3 5 0.3	∞ 5 ∞	5 ∞ ∞ 0.2	∞ ∞ ∞
$\rho_j(\alpha)$	{R2} {R1,R2}	0.2	0.2 0.15	0.3 0.2	0.3 0.25	0.2 0	0 0	0

*Note.* We consider high/low prices with and without bundling discount (i.e., buying {R1, R2} is cheaper than buying each individually). The other assortments can be used if items sell out.

where the first component is the remaining knapsack capacity. At a first stage, given a state of the form  $s = (b, \emptyset)$ , we choose action  $\diamond \in \{a, r\}$ , reducing the capacity if  $\diamond = a$ . At the second stage, the state is of the form  $s = (b, \diamond)$  with  $\diamond \in \{a, r\}$ , and we collect the reward only if  $\diamond$  = a. Formally, the rewards are  $\mathcal{R}((b, \mathbf{a}), \xi^{t-1/2}, \emptyset) = \xi^{t-1/2} \text{ and } \mathcal{R}((b, \mathbf{r}), \xi^{t-1/2}, \emptyset) = 0.$ 

# 6.1. Offline Benchmark and Online Policy for

To define Offline,  $\varphi$  and  $\hat{\varphi}$ , consider the following LP parametrized by  $(b, \mathbf{y}, \mathbf{z}) \in \mathbb{R}_{\geq 0} \times \mathbb{R}^n \times \mathbb{R}^n_{>0}$ :

6.1. Offline Benchmark and Online Policy for Distribution-Agnostic Online Knapsack

To define Offline, 
$$\varphi$$
 and  $\hat{\varphi}$ , consider the following LP parametrized by  $(b, \mathbf{y}, \mathbf{z}) \in \mathbb{R}_{\geq 0} \times \mathbb{R}^n \times \mathbb{R}_{\geq 0}^n$ :

$$(P[b, \mathbf{y}, \mathbf{z}]) \quad \text{maximize:} \quad \sum_j y_j x_j$$

$$\text{subject to:} \quad \sum_j w_j x_{ja} \leq b,$$

$$x_{ja} + x_{jr} = z_j \quad \forall j \in [n],$$

$$\mathbf{x} \geq 0.$$

Note that if the average rewards  $\mathbf{r}$  were known, then

Note that if the average rewards r were known, then setting y = r, we get the LP relaxation of Equation (1) for the baseline Online Knapsack. Moreover, for any r, the optimal LP solution sorts types by their "bang for the buck" ratios  $r_i/w_i$ , and accepts them greedily.

In particular, the solution only requires knowing the ranking induced by r.

#### 6.1.1. Offline Benchmark and Relaxed Value Function.

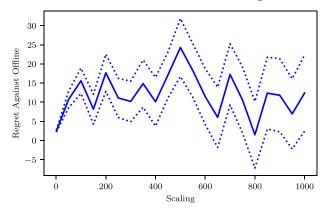
In this setting, we define Offline as the controller that knows the number of arrivals  $Z_i^T$  for each j, and also knows the ranking of the types (i.e., knows  $r_i/w_i \ \forall \ j \in [n]$ ).

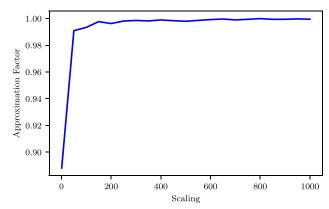
Formally, Offline is defined via the filtration  $G_t$  =  $\sigma(\{\xi^t:t\in[T]\}\cup\{\xi^\tau:\tau\geq t\})$ . This is a canonical filtration (see Definition 1) with variables  $(G_{\theta} : \theta \in \Theta) =$  $(\xi^t: t \in [T])$ . Observe that the future rewards, corresponding to times t - 1/2, are not revealed. Moreover, the relaxed value is defined as  $\varphi(t, s|\mathcal{G}_t) = P[b, \mathbf{r}, Z^t]$  for first stages and

$$\varphi(t-1/2,s|\mathcal{G}_t) = \begin{cases} P[b,\mathbf{r},Z^{t-1}] & \diamond = \mathbf{r} \\ \xi^{t-1/2} + P[b,\mathbf{r},Z^{t-1}] & \diamond = \mathbf{a}. \end{cases}$$
(16)

Remark 3 (MDP Relaxations for Distribution-Agnostic Settings). We note here that the underlying problem in this setting does not directly admit an MDP, as the distribution of rewards is unknown. However, once we reveal the arrivals to Offline, the relaxation does admit a well-defined MDP. By benchmarking against Offline, we bypass the need to explicitly formulate

Figure 7. (Color online) Performance of Pricing rabbi with Customer Choice (See Table 1)





Notes. We set the horizon as T = 10k and the inventory as (R1, R2) = (3k, 2k), and vary scaling parameter  $k = 1, \dots, 1,000$ . (Left) Regret against Offline, with 90% confidence intervals. (Right) Approximation ratio of RABBI against DP.

an Online control problem with distribution learning in this setting.

### 6.1.2. Value Function Estimate and Online Policy.

Recall we define  $R_j^t$  to be the empirical average of the observed rewards for type-j with t periods to go. We define the estimated value as  $\hat{\varphi} = P[B^t, R^t, \mathbb{E}[Z^t]]$ , resulting in the corresponding online policy given in Algorithm 4.

### Algorithm 4 (Learning RABBI)

**Input:** Access to solutions of (P[b, y, z])

Output: Sequence of decisions for Online.

- 1: Set  $B^T \leftarrow B$  as the given initial state and  $R^T$  as the single sample of each j.
- 2: **for**  $t \in \{T, T-1, ..., 1\}$  **do**
- 3: Compute  $X^t$ , an optimal solution to  $(P[B^t, R^t, \mathbb{E}[Z^t]])$ .
- 4: Observe the arrival type (context), say  $\xi^t = j$ , and take any action  $\hat{U}^t \in \arg\max_{u=a,r} \{X_{ju}^t\}$ .
- 5: If  $\hat{U}^t = a$ , collect random reward  $R_j$  and reduce the budget  $B^{t-1} \leftarrow B^t w_j$ . Else,  $B^{t-1} \leftarrow B^t$ .
- 6: Update empirical averages  $R^{t-1}$  based on  $R^t$  and the observation  $R_i$ .

# 6.2. Regret Analysis for Distribution-Agnostic Online Knapsack

As in the earlier sections, we first demonstrate that  $\phi$  satisfies the Bellman inequalities.

**Lemma 8.** The relaxation  $\varphi$  defined in (16) satisfies the Bellman inequalities with the following exclusion sets:

$$\mathcal{B}(t,b) = \{ \omega \in \Omega : \not\exists X \text{ solving } (P[b,\mathbf{r},Z^t])$$
s.t.  $X_{\mathcal{E}^t,\mathbf{a}} \ge 1 \text{ or } X_{\mathcal{E}^t,\mathbf{r}} \ge 1 \}.$ 

**Proof.** The initial ordering in Definition 2 follows from an argument identical to that of Lemma 2. The monotonicity property follows from Proposition EC.1 in the online appendix.  $\Box$ 

To complete the proof of Theorem 4, we need to characterize the information loss under Algorithm 4. The relaxation relies on the knowledge of  $\mathbf{r}$  (the true expectation) and  $Z^t$ . The natural estimators are the empirical averages  $R^t$  and expectation  $\mu(t) = \mathbb{E}[Z^t]$ , respectively. Specifically, we use maximizers  $X^t$  of  $(P[b, R^t, \mu(t)])$  to guess those of  $(P[b, \mathbf{r}, Z^t])$ .

The overall regret bound is  $r_{\varphi}(\text{Regret}_1 + \text{Regret}_2)$ , where  $\text{Regret}_1$  and  $\text{Regret}_2$  are two specific sources of error. When the estimators  $R^t$  of  $\mathbf{r}$  are accurate enough, the error is  $\text{Regret}_1$  and is attributed to the incorrect guess of a satisfying action, that is,  $\text{Regret}_1$  is an algorithmic regret. The second term,  $\text{Regret}_2$ , is the error that arises from insufficient accuracy of  $R^t$ , that is,  $\text{Regret}_2$  is the learning regret. The maximum

loss satisfies  $r_{\varphi} \leq \max_{j,i} \{w_i r_j / w_j - r_i\}$  and we can show that

Regret<sub>1</sub> 
$$\leq 2 \sum_{j} \frac{(w_{\text{max}}/w_{j})^{2}}{p_{j}}$$
 and Regret<sub>2</sub>  $\leq 16 \sum_{j} \frac{1}{p_{j}(w_{j}\delta)^{2}}$ .

In sum, the regret is bounded by  $(\max_{j,i} \{w_i r_j / w_j - r_i\}) \cdot (2 \sum_j \frac{(w_{\max}/w_j)^2}{p_i} + 16 \sum_j \frac{1}{p_i(w_i\delta)^2})$ .

**Remark 4** (Non-i.i.d Arrival Processes). We used the i.i.d. arrival structure to bound two quantities in the proof of Theorem 4: (1)  $\mathbb{P}[\|Z^t - \mathbb{E}[Z^t]\| \ge c\mathbb{E}[Z^t]]$ , and (2)  $\mathbb{E}[e^{-cN_j^t}]$ , where, recall,  $N_j^t$  is the number of type-j observations. The result holds for other arrival processes that admit these tail bounds.

#### 6.3. Censored Feedback

We consider now the case where only accepted arrivals reveal their reward. We retain the assumption of Theorem 4 that there is a separation  $\delta > 0$ :  $\bar{r}_j - \bar{r}_{j'} \ge \delta$  for all  $j \ne j'$ , where  $\bar{r}_i = \mathbb{E}[R_i]/w_i$ .

In the absence of full feedback, we will introduce a unified approach to obtaining the optimal regret (up to constant factors), that takes the learning method as a plug-in. The learning algorithm will decide between explore or exploit actions. Examples of learning algorithms that also give bounds that are explicit in t include modifications of UCB (Wu et al. 2015),  $\varepsilon$ -greedy or simply to set apart some time for exploration (see Corollary 1).

Recall that  $\sigma:[n] \to [n]$  is the ordering of [n] w.r.t. the ratios  $\bar{r}_j = r_j/w_j$  and  $\hat{\sigma}^t:[n] \to [n]$  is the ordering w.r.t. ratios  $\bar{R}_j^t = R_j^t/w_j$ . The discrepancy  $\mathbb{P}[\sigma \neq \hat{\sigma}^t]$  depends on the plug-in learning algorithm (henceforth Bandits). Bandits receives as inputs the current state  $S^t$  (remaining capacity), time, and the natural filtration  $\mathcal{F}_t$ . The output of Bandits is an action in {explore, exploit}. If the action is explore, we accept the current arrival in order to gather information, otherwise we call our algorithm to decide, as summarized in Algorithm 5. Note that  $\mathcal{F}_t$  has information only on the observed rewards, that is, accepted items.

### Algorithm 5. Bandits RABBI

**Input:** Access to Bandits and Algorithm 4.

Output: Sequence of decisions for Online.

- 1: Set  $S^T$  as the given initial state 2: **for** t = T, ..., 1 **do**
- 3: Observe input  $\xi^t$  and let  $U \leftarrow \text{BANDITS}(T, t, S^t, \mathcal{F}_t)$ .
- 4: If U = explore, accept the arrival.
- 5: If U = exploit, take the action given by Algorithm 4.
- 6: Update state  $S^{t-1} \leftarrow S^t w_{\xi^t}$  if accept or  $S^{t-1} \leftarrow S^t$  if reject.

**Theorem 7.** Let Regret<sub>1</sub> be the regret of Algorithm 4, as given in Theorem 4. Define the indicators explore<sub>t</sub>, exploit<sub>t</sub> which denote the output of Bandits at time t. The regret of Algorithm 5 is at most  $r_{\omega}M$ , where

$$M = \mathrm{Regret}_1 + \mathbb{E}\left[\sum_t explore_t\right] + \mathbb{E}\left[\sum_t \mathbb{P}\left[\sigma \neq \hat{\sigma}^t\right] exploit_t\right].$$

The expected regret of Algorithm 5 is thus bounded by the regret of Algorithm 4 in the full feedback setting, plus a quantity controlled by Bandits. In the periods where Bandits says explore (which, in particular, implies accepting the item), the decision might be the wrong one (i.e., different from Offline's). We upper bound this by the number of exploration periods. This is the second term in M. The decision might also be wrong if Bandits says exploit (in which case we call Algorithm 4), but the (learned) ranking at time t,  $\hat{\sigma}^t$ , is different from  $\sigma^t$ . This is the last term in M. Finally, even if the learned ranking is correct, exploit can lead to the wrong guess by Algorithm 4 because the arrival process is uncertain. This is the first term in M.

Corollary 1 uses a naive Bandits, which explores until obtaining  $\Omega(\log T)$  samples and achieves the optimal (i.e., logarithmic) regret scaling. The constants may be improved by changing the Bandits module we use. Any such algorithm has the guarantee given by Theorem 7. With the naive Bandits, the bound follows from a generalization of coupon collector (Shank and Yang 2013).

**Corollary 1.** If we first obtain  $\frac{8}{(w_i\delta)^2}\log T$  samples of every type j, then we can obtain  $O(\log T)$  regret, which is optimal up to constant factors.

### 7. Concluding Remarks

We developed a framework that provides rigorous support to the use of simple optimization problems as a basis for online resolving algorithms. The framework is based on using a carefully chosen offline benchmark, which guides the online algorithm. The regret bounds then follow from our use of Bellman inequalities and a useful distinction between Bellman loss and information loss.

As is often the case in approximate dynamic programming, the identification of a function  $\varphi$  satisfying the Bellman inequalities requires some ad hoc creativity but, as our examples illustrate, is often rather intuitive. In online Appendix EC.1, we provide sufficient conditions, applicable to cases where  $\varphi$  has a natural linear representation, to verify the Bellman inequalities. These conditions are intuitive and likely to hold for a variety of resource-allocation problems. Importantly, once such a function is identified, our RABBI framework provides a way of obtaining

online policies from  $\varphi$ , and corresponding regret bounds.

We illustrate our framework on three settings. First, we consider online probing, which serves as an instance of a larger family of two-stage decision problems, wherein there is an inherent trade-off between getting refined information, and the cost of obtaining it. Next, we consider dynamic pricing, which is a well-studied problem, and is representative of settings where rewards and transitions are random. Finally, our study of online contextual bandits with knapsacks showcases a separation of the underlying combinatorial problem from the parameter estimation problem.

It is our hope that this structured framework will be useful in developing online algorithms for other problems, whether these are extensions of those we studied here or completely different.

### References

Agrawal S, Devanur N (2016) Linear contextual bandits with knapsacks. *Proc. 30th Internat. Conf. Neural Inform. Processing Systems* (Curran Associates Inc, Red Hook, NY), 3450–3458.

Alaei S (2014) Bayesian combinatorial auctions: Expanding single buyer mechanisms to many buyers. SIAM J. Comput. 43(2):930–972.

Arlotto A, Gurvich I (2019) Uniformly bounded regret in the multisecretary problem. Stochastic Systems 9(3):231–260.

Babaioff M, Dughmi S, Kleinberg R, Slivkins A (2015) Dynamic pricing with limited supply. *ACM Trans. Econom. Comput.* 3(1): article 4, https://doi.org/10.1145/2559152.

Badanidiyuru A, Kleinberg R, Slivkins A (2018) Bandits with knapsacks. *J. ACM* 65(3):Article 13.

Badanidiyuru A, Langford J, Slivkins A (2014) Resourceful contextual bandits. Conf. Learn. Theory, 1109–1134.

Balseiro SR, Brown DB (2019) Approximations to stochastic dynamic programs via information relaxation duality. Oper. Res. 67(2): 577–597.

Banerjee S, Freund D (2020) Uniform loss algorithms for online stochastic decision-making with applications to bin packing. *Abstr. 2020 SIGMETRICS/Performance Joint Internat. Conf. Mea surement Model. Comput. Systems* (Association for Computing Machinery, New York), 1–2.

Boucheron S, Lugosi G, Massart P (2013) Concentration Inequalities: A Nonasymptotic Theory of Independence (Oxford University Press, Oxford, UK).

Brown D, Smith J, Sun P (2010) Information relaxations and duality in stochastic dynamic programs. *Oper. Res.* 58(4):785–801.

Bubeck S, Cesa-Bianchi N (2012) Regret analysis of stochastic and nonstochastic multi-armed bandit problems. Foundations Trends Machine Learn. 5(1):1–122.

Bubeck S, Perchet V, Rigollet P (2013) Bounded regret in stochastic multi-armed bandits. *Proc. Conf. Learn. Theory*, PLMR 30: 122–134.

Buchbinder N, Jain K, Singh M (2014) Secretary problems via linear programming. *Math. Oper. Res.* 39(1):190–206.

Bumpensanti P, Wang H (2020) A re-solving heuristic for dynamic resource allocation with uniformly bounded revenue loss. *Management Sci.*, ePub ahead of print March 16, https://doi.org/ 10.1287/mnsc.2019.3365.

Chen Q, Jasin S, Duenyas I (2019) Nonparametric self-adjusting control for joint learning and optimization of multiproduct pricing with finite resource capacity. Math. Oper. Res. 44(2): 377–766.

- Chugg B, Maehara T (2019) Submodular stochastic probing with prices. *Proc. 66th Internat. Conf. Control, Decision Inform. Tech.* (IEEE, Paris), 60–66.
- Correa J, Foncea P, Hoeksma R, Oosterwijk T, Vredeveld T (2017) Posted price mechanisms for a random stream of customers. *Proc. ACM Conf. Econom. Comput.*, 169–186.
- Düetting P, Feldman M, Kesselheim T, Lucier B (2017) Prophet inequalities made easy: Stochastic optimization by pricing non-stochastic inputs. *Proc. 58th Annual Symp. Foundations Comput. Sci.* (IEEE, Berkeley, CA), 540–551.
- Gallego G, Van Ryzin G (1997) A multiproduct dynamic pricing problem and its applications to network yield management. *Oper. Res.* 45(1):24–41.
- Gupta A, Nagarajan V (2013) A stochastic probing problem with applications. Proc. Internat. Conf. Integer Programming Combin. Optim. (Springer, Berlin, Heidelberg), 205–216.
- Gupta A, Nagarajan V, Singla S (2016) Algorithms and adaptivity gaps for stochastic probing. Proc. 27th Annual ACM-SIAM Symp. Discrete Algorithms, 1731–1747.
- Jasin S (2014) Reoptimization and self-adjusting price control for network revenue management. *Oper. Res.* 62(5):1168–1178.
- Jasin S, Kumar S (2012) A re-solving heuristic with bounded revenue loss for network revenue management with customer choice. Math. Oper. Res. 37(2):313–345.
- Kleinberg R, Weinberg SM (2012) Matroid prophet inequalities. *Proc.* 44th Annual ACM Symp. Theory Comput., 123–136.
- Mangasarian O, Shiau T (1987) Lipschitz continuity of solutions of linear inequalities, programs and complementarity problems. *SIAM J. Control Optim.* 25(3):583–595.
- Massart P (1990) The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *Ann. Probab.* 18(3):1269–1283.

- Shank N, Yang H (2013) Coupon collector problem for non-uniform coupons and random quotas. *Electron. J. Combin.* 20(2), doi:.
- Singla S (2018) The price of information in combinatorial optimization. Proc. 29th Annual ACM-SIAM Symp. Discrete Algorithms, 2523–2532.
- Talluri KT, Van Ryzin GJ (2006) The Theory and Practice of Revenue Management, vol. 68 (Springer Science & Business Media, New York).
- Vera A, Banerjee S (2020) The Bayesian prophet: A low-regret framework for online decision making. *Management Sci.*, ePub ahead of print October 5, https://doi.org/10.1287/mnsc.2020.3624.
- Weitzman ML (1979) Optimal search for the best alternative. *Econometrica* 47(3):641–654.
- Wu H, Srikant R, Liu X, Jiang C (2015) Algorithms with logarithmic or sublinear regret for constrained contextual bandits. Adv. Neural Inform. Processing Systems 28:433–441.

**Alberto Vera** received his PhD in Operations Research from Cornell in 2020. He is currently a research scientist at Amazon working on resource allocation problems.

**Sid Banerjee** is an assistant professor in the School of Operations Research at Cornell, working on topics at the intersection of data-driven decision-making, network algorithms and market design. He received his PhD in 2013 from the ECE Department at UT Austin, and was a postdoctoral researcher in the Social Algorithms Lab at Stanford from 2013-2015

**Itai Gurvich** is a professor at Northwestern's Kellogg School of Management. His research focuses on the performance analysis and optimization of processing networks.