

Online Allocation and Pricing: Constant Regret via Bellman Inequalities

Alberto Vera Siddhartha Banerjee Itai Gurvich
School of Operations Research and Information Engineering, Cornell University
aav39@cornell.edu sbanerjee@cornell.edu gurvich@cornell.edu

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 re-solving an LP in each period, followed by a simple greedy action-selection rule; thus, our policies are practical as well as provably near optimal.

Key words: 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 (MDP), 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, etc.

The above 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, i.e., 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 requirement-reward 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 `OnlineKnapsack` 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 to the benchmark (and thus, compared to 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 also, 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 which 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 who 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 which 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 Overview

2.1. Problem Settings and Results

We illustrate our framework by developing low-regret 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 $t = T, T - 1, \dots, 1$ (where t denotes the *time-to-go*), we assume the arriving type is drawn from a finite set $[n]$ from some known distribution $\mathbf{p} = (p_1, \dots, p_n)$. 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 , the controller has an additional probing budget B_p 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 $\{f_{j1}, \dots, f_{jm}\}$; 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 (*Re-solve and Act Based on the Bellman Inequalities*; see Section 3.2) is based on comparing two ‘controllers’: OFFLINE, who acts optimally given future information, and a non-anticipative controller ONLINE who tries to follow OFFLINE. Both start in the same initial state S^T . We denote v^{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 non-anticipative policy π that maps current states to actions, resulting in a total expected reward v_π^{on} .

Let π_R denote the online policy produced by our RABBI framework, and π denote any non-anticipative policy. Then the expected regret of π_R relative to the chosen offline benchmark is

$$\mathbb{E}[\text{Regret}] := v^{\text{off}} - v_{\pi_R}^{\text{on}} \geq \max_{\pi} [v_\pi^{\text{on}}] - v_{\pi_R}^{\text{on}}$$

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

For all the above 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), 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 .*

The above 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 $(A, \mathbf{f}, F_1, \dots, F_n)$, 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 Fig. 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

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

The function $V(\cdot|Z^t)$ is thus `OFFLINE`'s value function (see Fig. 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.

$V(\cdot|Z^t)$ can also be represented via Bellman equations. Specifically, at time-to-go t , assuming `OFFLINE` has budget b and the arriving type is ξ , the value function obeys the Bellman equation

$$V(t, b|Z^t) = \max \left\{ [r_{\xi t} + V(t-1, b - w_{\xi t}|Z^{t-1})] \mathbb{1}_{\{w_{\xi t} \leq b\}}, V(t-1, b|Z^{t-1}) \right\}, \quad \forall t, b, \xi^t.$$

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

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

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

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

Note that, if $\mathbb{E}[L_B(t, b)]$ is small, with expectation taken over Z^t , then φ ‘almost’ satisfies the Bellman equations. We henceforth refer to φ as a *relaxed value* for V and L_B the *Bellman Loss*.

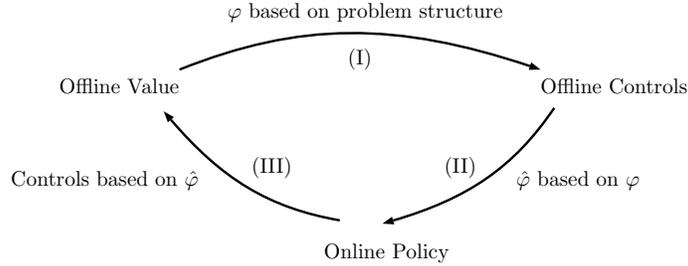


Figure 1 The RABBI framework: We first define OFFLINE’s value function by specifying access to future information. Next, we identify a tractable relaxation φ for OFFLINE’s value under this same information structure (step I). Finally, we introduce a non-anticipative estimate $\hat{\varphi}$ for φ , and use it to design online controls (step II). The resulting online policy is evaluated against OFFLINE’s value (step III).

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

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

Note that $\hat{\varphi}$ does not approximate V or φ 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 $\text{argmax}_{u \in \{a, r\}} \{y_{\xi^t, u}\}$) guarantees that that ONLINE and OFFLINE play the same action with high probability. Whenever OFFLINE and ONLINE play different actions, then 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 Fig. 1.

Towards 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 φ 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 *Re-solve 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, Vera and Banerjee 2020, Banerjee and Freund 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 to the full information benchmark* (see Example 1).

- **Explicit value-function characterizations:** The optimization problem in Eq. (1) has a closed-form solution, which was used explicitly by (Arlotto and Gurvich 2019, Vera and Banerjee 2020, Banerjee and Freund 2020); this does not extend to more complex settings.

Our framework in this work resolves these shortcomings 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, Düetting et al. 2017, Correa 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 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, 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 Appendix E.

Moreover, the different problems we apply RABBI each have a large body of prior work.

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 *non-degenerate* instances. This inefficiency, though, is fundamental, since they use the ex ante (or fluid) benchmark, which has $\Omega(\sqrt{T})$ under non-degeneracy. 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 i.i.d. Poisson arrivals, and require 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 QOS constraints. See Vera and Banerjee (2020), Banerjee and Freund (2020) for more discussion and references.

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 *non-adaptive* constant-factor competitive algorithms for this problem (Gupta et al. 2016). In terms of *online adaptive* algorithms, Chugg and Maehara (2019) introduces an algorithm with bounded competitive ratio in an adversarial setting.

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) provides a policy with $O((B \log T)^{2/3})$ regret under *adversarial inputs*, as opposed to our $O(1)$ guarantee under stochastic inputs.

Knapsack with learning. Multi-armed bandit problems have been widely studied, and we refer to Bubeck et al. (2013, 2012) 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 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 $t = T, T-1, \dots, 1$; 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 ξ . This assumption is for ease of exposition only; our results can be extended to hold when rewards or transitions are random given ξ .

The feasible actions for state s and input ξ correspond to the set $\{u \in \mathcal{U} : \mathcal{R}(s, \xi, u) > -\infty\}$. We assume that this feasible set is non-empty for all $s \in \mathcal{S}, \xi \in \Xi$, and also, that the maximum reward is bounded, i.e., $\sup_{s \in \mathcal{S}, \xi \in \Xi, u \in \mathcal{U}} \mathcal{R}(s, \xi, u) < \infty$.

The MDP described above induces a natural filtration \mathcal{F} , with $\mathcal{F}_t = \sigma(\{\xi^\tau : \tau \geq t\})$; a non-anticipative policy is one which 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^{off} is an upper bound on the performance of the optimal non-anticipative policy.

We present a specific class of filtration (generated by augmenting the canonical filtration) that suffice for our applications (see Fig. 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_Θ 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 , i.e., 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 Eq. (2) reduces to the Bellman equation for the value-function of an optimal non-anticipative policy:

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

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

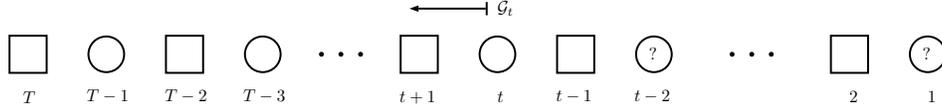


Figure 2 Illustration of Definition 1. In online probing (see Section 4), arrivals first reveal their public type, then the controller chooses an action (accept/probe/reject), and then the private type (true reward) is revealed. Squares (resp. circles) represent public (resp. private) information. The filtration \mathcal{G} used by RABBI comprises of all public types, i.e. $G_\theta = (\xi^\theta : \xi^\theta \text{ is a public type})$. At time t , OFFLINE knows all the information thus far (to the left and including t), plus the future squares.

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 φ and Bellman Inequalities. Intuitively, φ is “almost” defined by a dynamic-programming recursion; quantitatively, whenever φ 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 r.v. $\{\varphi(t, s)\}_{t,s}$ satisfies the Bellman Inequalities w.r.t. filtration \mathcal{G} and r.v. $\{L_B(t, s)\}_{t,s}$ 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: $\forall s \in \mathcal{S}, t \in [T]$,

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

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

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

A special case is when the Bellman Loss is 0 over sample paths in some chosen set:

DEFINITION 3 (EXCLUSION SETS). A set $\mathcal{B}(t, s)$ is an *exclusion set* if we can write the Bellman Loss as $L_B(t, s) = r_\varphi \mathbf{1}_{\mathcal{B}(t, s)}$ 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 $\varphi(t, s|\mathcal{G}_t) \leq \max_{u \in \mathcal{U}} \{\mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u)|\mathcal{G}_{t-1})|\mathcal{G}_t]\}$, i.e., 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 $(P[\mathbf{d}]) : \max\{\mathbf{r}'\mathbf{x} : M\mathbf{x} = \mathbf{d}, \mathbf{x} \geq 0\}$, where $M \in \mathbb{R}^{m \times n}$ is an arbitrary constraint matrix. If $\bar{\mathbf{x}}$ solves $(P[\mathbf{d}])$ and $\bar{x}_j \geq 1$ 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 $x_j \geq 1$. Therefore we have $P[\mathbf{d}] = \max\{\mathbf{r}'(\mathbf{x} + \mathbf{e}_j) : M(\mathbf{x} + \mathbf{e}_j) = \mathbf{d}, \mathbf{x} \geq 0\}$. \square

Lemma 1 lets us divide $P[\mathbf{d}]$ in 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

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

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

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

Moreover, for our choice of φ , 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 $L_B(t, b) = r_\varphi \mathbf{1}_{\mathcal{B}(t, b)}$. \blacksquare

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

DEFINITION 4 (MAXIMUM LOSS). For a given relaxation φ , the maximum loss is given by

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

Next, note that the ‘optimal’ action in the RHS of Eq. (3) need not be unique, and indeed the inequality can be satisfied by multiple actions. For given φ and L_B , we define:

DEFINITION 5 (SATISFYING ACTIONS). Given a filtration \mathcal{G} and relaxed value φ , we say that u is a *satisfying action* for state s at time t if

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

At any time t and state $s \in \mathcal{S}$, any action in $\operatorname{argmax}_{u \in \mathcal{U}} \{\mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u) | \mathcal{G}_{t-1}) | \mathcal{G}_t]\}$ 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 φ 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 the above 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 $\mathcal{G}_t = \sigma(\{\xi^l : l \geq t\} \cup G_\Theta)$, and assume that the relaxed value φ can be represented as a function of the random variables $\{\xi^l : l \geq t\} \cup G_\Theta$ 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 φ is expressed as the solution of an optimization problem:

$$\varphi(t, s; f_t(\xi^T, \dots, \xi^t, G_\Theta)) = \max_{\mathbf{x} \in \mathbb{R}^{\mathcal{U} \times \Xi}} \{h_t(\mathbf{x}; s, f_t(\xi^T, \dots, \xi^t, G_\Theta)) : g_t(\mathbf{x}; s, f_t(\xi^T, \dots, \xi^t, G_\Theta)) \leq 0\}. \quad (5)$$

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 ξ in remaining periods. We can also interpret $x_{u,\xi}$ as a *score* for action u when input ξ is presented. Now to get a non-anticipative policy, a natural ‘projection’ of $\varphi(t, s | \mathcal{G}_t)$ on the filtration \mathcal{F} is given via the following optimization problem

$$\hat{\varphi}(t, s | \mathcal{F}_t) = \varphi(t, s; \mathbb{E}[f_t(\xi^T, \dots, \xi^t, G_\Theta) | \mathcal{F}_t]) = \max_{\mathbf{y} \in \mathbb{R}^{\mathcal{U} \times \Xi}} \{h_t(\mathbf{y}; s, \mathbb{E}[f_t | \mathcal{F}_t]) : g_t(\mathbf{y}; s, \mathbb{E}[f_t | \mathcal{F}_t]) \leq 0\}. \quad (6)$$

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

RABBI (Re-solve and Act Based on Bellman Inequalities)

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

Output: Sequence of decisions \hat{U}^t for ONLINE.

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

$$\mathbb{E}[\text{Regret}] \leq \mathbb{E} \left[\sum_t (r_\varphi \mathbb{1}_{\mathcal{Q}(t, S^t)} + \mathbb{1}_{\mathcal{Q}(t, S^t)^c} L_B(t, S^t)) \right] \leq \sum_t (r_\varphi \mathbb{P}[\mathcal{Q}(t, S^t)] + \mathbb{E}[L_B(t, S^t)]).$$

REMARK 1 (BELLMAN AND INFORMATION LOSS). The bound in Theorem 5 has two distinct summands: The *information loss* $\sum_t \mathbb{P}[\mathcal{Q}(t, S^t)]$ measures how often RABBI takes a non-satisfying action due to randomness in sample paths; on the other hand, the *Bellman loss* $\sum_t \mathbb{E}[L_B(t, S^t)]$ quantifies violations of the Bellman equations made under the pseudo value-function φ . ■

Compensated Coupling: 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’s always sufficient to compensate OFFLINE by the *maximum loss* r_φ 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 of three different types, and the sequence of actions $(\mathbf{r}, \mathbf{a}, \mathbf{r}, \mathbf{a}, \mathbf{r})$ ((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 collected 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 $\forall t$:

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

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:

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

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_φ . Specifically, we have

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

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

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

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 $R_j \in \{r_{jk} : k \in [m]\}$ drawn with probabilities $\{q_{jk}\}$; \mathbf{r} and \mathbf{q} are known. We assume w.l.o.g that $r_{j1} < r_{j2} < \dots < r_{jm}$ and $r_{jm} > 0$. For ease of exposition, we assume that all arrivals have unit weights; our analysis however extends to general weights w_j . The controller may accept (a), reject (r) or probe (p) the arrival. Accepting type- j item without probing results in expected reward of $\bar{r}_j := \sum_{k \in [m]} r_{jk} q_{jk}$. 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 (resp. probed), we reduce B_h (resp. B_p) by one.

Formally, we view each time period $t \in \{T, T-1, \dots, 1\}$ 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 $\{\mathbf{a}, \mathbf{p}, \mathbf{r}\}$;

in the second stage $t - 1/2$, the reward r_{jk} (or “sub-type” $\xi^{t-1/2} = (j, k) \in [n] \times [m]$) is drawn with probability q_{jk} , and the available actions are $\{\mathbf{a}, \mathbf{r}\}$ if the first-stage action is \mathbf{p} , and \emptyset otherwise. We augment the state space with a variable \diamond that captures the first stage decision (i.e., whether we accept/reject without probing or probe). The state space \mathcal{S} of the controlled process is thus $\mathcal{S} = \{(b_h, b_p, \diamond) : b_h, b_p \in \mathbb{N}, \diamond \in \{\mathbf{a}, \mathbf{p}, \mathbf{r}, \emptyset\}\}$, where b_h, b_p are the residual hiring and probing budgets. In first stage of each period, we set $\diamond = \emptyset$, and only collect rewards in second stage in each period. See Fig. 3 for an illustration.

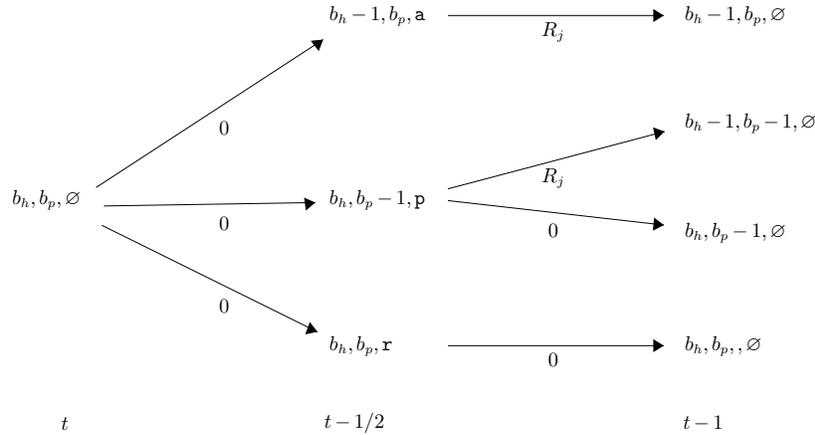


Figure 3 Actions/transitions in online probing in periods t , $t - 1/2$, and $t - 1$, with inputs $\xi^t = j$ and $\xi^{t-1/2} = R^j$. Numbers below the arrows represent the reward of a transition. At t , available actions are $\{\mathbf{a}, \mathbf{p}, \mathbf{r}\}$ (i.e., accept, probe, reject; from 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 $\{\mathbf{a}, \mathbf{r}\}$.

4.1. Offline Benchmark and Online Policy for Probing

We now apply the RABBI framework for online probing.

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 (sub-types). 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 or not to probe an arrival.

Relaxed Value Function: Since solving for OFFLINE's optimal actions may be non-trivial, we next construct a relaxed value function φ , using the following LP parametrized by $(b_h, b_p, \mathbf{z}) \in \mathbb{N}^2 \times \mathbb{R}_{\geq 0}^n$,

$$\begin{aligned}
(P[b_h, b_p, \mathbf{z}]) \quad & \text{maximize:} && \sum_{j,k} r_{jk} x_{jka} + \sum_j \bar{r}_j x_{ja} && (7) \\
& \text{subject to:} && \sum_{j,k} x_{jka} + \sum_j x_{ja} \leq b_h \\
& && \sum_j x_{jp} \leq b_p \\
& && x_{ja} + x_{jp} + x_{jr} = z_j \quad \forall j \in [n] \\
& && x_{jka} + x_{jkr} = q_{jk} x_{jp} \quad \forall j \in [n], k \in [m] \\
& && \mathbf{x} \geq 0
\end{aligned}$$

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 $\mathbf{x} \in \mathbb{R}_{\geq 0}^{3n+2nm}$ represent action summaries, where x_{ja}, x_{jr}, 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 sub-type k (i.e., reward r_{jk}).

To construct relaxed value φ , recall that a state is of the form $s = (b_h, b_p, \diamond)$ with $\diamond \in \{\mathbf{a}, \mathbf{p}, \mathbf{r}, \emptyset\}$. 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 φ to incorporate the action $(\mathbf{a}, \mathbf{p}, \mathbf{r})$ taken in the first stage. Overall, our relaxation is defined as follows

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

Value Function estimate and Online Policy: Finally we can use the relaxed value function ϕ in Eq. (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.

REMARK 2 (PROBING COST). Our approach can also handle a setting where the controller has no probing budget, but instead incurs a penalty c_j when probing a type- j arrival. 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 $\max\{\sum_{j,k} r_{jk} x_{jka} + \sum_j \bar{r}_j x_{ja} - \sum_j c_j x_{jp}\}$ \blacksquare

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

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 Appendix B.

The main part of the proof involves showing that φ as defined in Eq. (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 φ .

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 φ . The proof is provided in Appendix B.

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 to 5).

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

- (1) *If $\bar{X}_{i\mathbf{a}} \geq 1$, then accepting at time t is a satisfying action.*
- (2) *If $\bar{X}_{i\mathbf{r}} \geq 1$, then rejecting at time t is a satisfying action.*
- (3) *If $\bar{X}_{i\mathbf{p}} \geq 1$, and $\xi^{t-1/2} = (i, k)$ is such that either $\bar{X}_{i\mathbf{k}\mathbf{a}} \geq 1$ or $\bar{X}_{i\mathbf{k}\mathbf{r}} \geq 1$, then probing at time t , followed by accepting (if $\bar{X}_{i\mathbf{k}\mathbf{a}} \geq 1$) or rejecting (if $\bar{X}_{i\mathbf{k}\mathbf{r}} \geq 1$) at time $t - 1/2$ is a satisfying action.*

Finally φ satisfies the Bellman Inequalities with Bellman Loss $L_B(t, (b_h, b_p)) = r_\varphi \mathbf{1}_{\mathcal{B}(t, b_h, b_p)}$, where \mathcal{B} are exclusion sets defined as:

$$\mathcal{B}(t, b_h, b_p) = \{\omega \in \Omega : \bar{X} \text{ solution to } (P[(b_h, b_p), Z^t]) \text{ s.t. (1) or (2) or (3) hold}\}.$$

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

$$\varphi(t, (b_h, b_p, \emptyset) | \mathcal{G}_t) \leq \max_{\diamond \in \{\mathbf{a}, \mathbf{p}, \mathbf{r}\}} \{\mathbb{E}_{\xi^{t-1/2}}[\varphi(t-1/2, (s_\diamond, \diamond) | \mathcal{G}_{t-1/2}) | \mathcal{G}_t]\} \quad \forall \omega \notin \mathcal{B}(t, b_h, b_p).$$

where the state $s_\diamond = (b_h - 1, b_p)$ if $\diamond = \mathbf{a}$, $s_\diamond = (b_h, b_p - 1)$ if $\diamond = \mathbf{p}$ and $s_\diamond = (b_h, b_p)$ if $\diamond = \mathbf{r}$. Moreover, given $\xi^t = i$, we have from Eq. (8) that $\mathbb{E}_{\xi^{t-1/2}}[\varphi(t-1/2, (s_\diamond, \diamond) | \mathcal{G}_{t-1/2}) | \mathcal{G}_t] = P[(b_h, b_p), Z^{t-1}]$ if $\diamond = \mathbf{r}$, and $r_{\xi^{t-1/2}} + P[(b_h - 1, b_p), Z^{t-1}]$ if $\diamond = \mathbf{a}$. Now for cases (1) and (2), the claim in the lemma follows directly by invoking Lemma 1. Finally, case (3) (where $\bar{X}_{\mathbf{ip}} \geq 1$) also follows from using Lemma 1, but in a somewhat more technical way; see Appendix B 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 $\text{Regret} \leq r_\varphi \sum_t (\mathbb{1}_{\mathcal{B}(t, S^t)} + \mathbb{1}_{\mathcal{Q}(t, S^t)})$. 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 \bar{X} be the solution to $(P[\mathbf{b}, Z^t])$, and note that Lemma 3 guarantees that there is zero Bellman Loss if (1) $\max\{\bar{X}_{j\mathbf{a}}, \bar{X}_{j\mathbf{r}}\} \geq 1$, or (2) $\bar{X}_{j\mathbf{p}} \geq 1$ and $\max\{\bar{X}_{j\mathbf{ka}}, \bar{X}_{j\mathbf{kr}}\} \geq 1$. 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 $x_{j\mathbf{a}} + x_{j\mathbf{p}} + x_{j\mathbf{r}} = Z_j^t \forall j$ and $x_{j\mathbf{ka}} + x_{j\mathbf{kr}} = q_{jk}x_{j\mathbf{p}} \forall j, k$. If $Z_j^t \geq 3$, then one of the variables $x_{j\mathbf{a}}, x_{j\mathbf{p}}, x_{j\mathbf{r}}$ must be at least 1. On the other hand, we need $q_{jk}x_{j\mathbf{p}} \geq 2$ to guarantee that one of $x_{j\mathbf{ka}}, x_{j\mathbf{kr}}$ is at least 1. Thus we have

$$\mathbb{P}[\mathcal{B}(t, b) | \xi^{t-1/2} = (j, k)] \leq \mathbb{P}\left[Z_j^t < \frac{6}{q_{jk}}\right] = \mathbb{P}\left[Z_j^t - \mu_j(t) < -\mu_j(t) \left(1 - \frac{6}{\mu_j(t)q_{jk}}\right)\right]. \quad (9)$$

Restricting $\mu_j(t) \geq 12/q_{jk}$ to ensure the RHS of Eq. (9) is positive, we can use a standard Chernoff bound (see (Boucheron et al. 2013)) to get $\mathbb{P}[\mathcal{B}(t, b) | \xi^{t-1/2} = (j, k)] \leq e^{-2(p_j/2)t} + \mathbb{1}_{\{t \leq 12/(p_j q_{jk})\}}$. 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 $\sum_t \mathbb{P}[\mathcal{Q}(t, S^t)]$, 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 $X_{j\mathbf{a}}^t + X_{j\mathbf{p}}^t + X_{j\mathbf{r}}^t = \mu_j(t)$, we have

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

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

$$\mathbb{P}\left[\bar{X}_{j\mathbf{p}} < 1 \text{ or } \bar{X}_{\xi^{t-1/2}, u} < 1 \mid X_{j\mathbf{p}}^t \geq \frac{\mu_j(t)}{3}, X_{\xi^{t-1/2}, u}^t \geq \frac{q_{\xi^{t-1/2}} \mu_j(t)}{6}\right] \leq \mathbb{P}\left[\|\bar{X} - X^t\|_\infty \geq \frac{q_{\xi^{t-1/2}} \mu_j(t)}{6}\right],$$

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 non-satisfying action is bounded by $\mathbb{P}[\|\bar{X} - X^t\|_\infty \geq \min_k q_{jk} \cdot \mu_j(t)/6]$. Moreover, standard LP sensitivity results (Mangasarian and Shiau 1987, Theorem 2.4) imply that there exists κ depending on \mathbf{q}, n, m alone, s.t. $\|\bar{X} - X^t\|_\infty \leq \kappa \|Z^t - \mu(t)\|_1$. Finally, the measure of sets \mathcal{Q} where ONLINE chooses a non-satisfying action is bounded by

$$\sum_t \mathbb{P}[\mathcal{Q}(t, S^t)] \leq \sum_t \mathbb{P}[\|Z^t - \mu(t)\|_1 \geq \min_k q_{jk} \cdot \mu_j(t)/6\kappa] < \infty.$$

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

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 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 make a choice 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 $\{f_{j1}, \dots, f_{jm}\}$. 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

Offline Benchmark: Note that for each customer type j , there are Z_j^T arrivals, and hence Z_j^T draws from the distribution F_j . We now define our benchmark by considering OFFLINE to be a controller that knows the realized histogram of these draws, i.e., for each j , OFFLINE knows the empirical distribution of the Z_j^T rewards. Moreover, at the end of each period t , OFFLINE also observes the realized valuation R^t whether or not there is a sale. Note that OFFLINE does not know the exact sequence of these rewards, and so is not a full information benchmark. For example, say

$Z_1^T = 15$ and we reveal that 10 arrivals type-1 have private reward \$1 and 5 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 $f_{j1} > f_{j2} > \dots > f_{jm}$. 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 $\{Y^t : t = T, T-1, \dots, 1\}$ where $Y_{jl}^t := \mathbb{1}_{\{\xi^t=j, R^t > f_{jl}\}}$; in other words, Y_{jl}^t is the indicator of whether a price f_{jl} or lower is accepted by the type- j at time t . We define $Q_{jl}(t) := \frac{1}{Z_j^t} \sum_{\tau=1}^t Y_{jl}^\tau$ 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 $Q_{jl}(t) = \frac{Z_j^{t+1}}{Z_j^t} Q_{jl}(t+1) - \frac{1}{Z_j^t} Y_{jl}^{t+1}$.

OFFLINE's information is now given by the filtration $\mathcal{G}_t = \sigma(\{Q(\tau), Z^\tau : \tau \geq t\})$, i.e., 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}, \dots, Y^1 instead of their empirical averages.

Relaxed Value Function: Consider the following LP, parameterized by $(\mathbf{b}, \mathbf{q}, \mathbf{z})$.

$$\begin{aligned}
 (P[\mathbf{b}, \mathbf{q}, \mathbf{z}]) \quad & \text{maximize:} && \sum_{j,l} f_{jl} q_{jl} x_{jl} && (10) \\
 & \text{subject to:} && \sum_{j,l} a_{ij} q_{jl} x_{jl} \leq b_i && \forall i \in [d] \\
 & && \sum_{j,l} x_{jl} + x_{jr} = z_j && \forall j \in [n] \\
 & && \mathbf{x} \geq 0 &&
 \end{aligned}$$

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 $q_{jl} \leftarrow \bar{F}_j(f_{jl})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: If the arrival is type j and $A_j \not\leq B^t$: 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 $l \in \operatorname{argmax}\{X_{j,l}^t : l = 1, \dots, m, \mathbf{r}\}$. If $l = \mathbf{r}$, 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 $q_{jl} = \bar{F}_j(f_{jl})$, i.e., 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 $\sum_l x_{jl} + x_{j\mathbf{r}} = z_j$, where $x_{j\mathbf{r}}$ is the number of rejected type- j customers.

5.2. Bellman Inequalities and Bellman Loss

We first argue that our choice of φ 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. $\mathbb{E}[V(T, B|\mathcal{G}_T)] \leq \mathbb{E}[\varphi(T, B|\mathcal{G}_T)]$, hence φ satisfies the initial ordering condition.
2. If the arriving type is j and $\max_l \{X_{jl}\} \geq 1$, then $\mathbb{E}[L_B(t, \mathbf{b})] \leq 0$.
3. If the arriving type is j and $X_{jl} \geq 1$, 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. Below we present the main ingredients for obtaining the monotonicity property (items (2) and (3)); complete details are deferred to Appendix C. For ease of exposition, when the controller rejects, he can equivalently post $f_{j\mathbf{r}} = \infty$ such that $\bar{F}_j(f_{j\mathbf{r}}) = 0$ 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 $\mathbf{b} - A_j Y_{jl}^t$, thus monotonicity corresponds to:

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

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

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

Define $L_B(t+1, \mathbf{b}, j, l) := \varphi(t+1, \mathbf{b}) - f_{jl} Q_{jl}(t+1) - \mathbb{E}_{t+1}[\varphi(t, \mathbf{b} - A_j Y_{jl}^{t+1})]$, which corresponds to the loss in Eq. (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 $P[\mathbf{b}, Q(t+1), Z^{t+1}] = f_{jl} Q_{jl}(t+1) + P[\mathbf{b} - A_j Q_{jl}(t+1), Q(t+1), Z^t]$. Thus, assuming $X_{jl} \geq 1$, we can write the loss in the Bellman inequality as

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

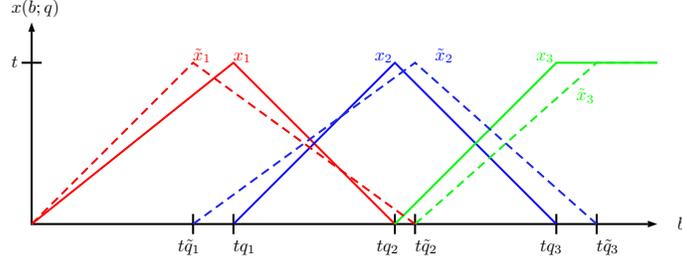


Figure 4 Solution to the pricing LP in Eq. (10) for the case $d = 1$ and $n = 1$, which correspond to selling multiple copies of an item to homogeneous customers. If $b/t \in (q_l, q_{l+1}]$, the prices used by the LP are f_l, f_{l+1} and the amount of time we offer each is piece-wise linear in the budget. 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} .

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 Eq. (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 Appendix C 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 $X_{jl} \geq 1$, $\mathbb{E}_{t+1}[L_B(t+1, b, j, l)] \leq 0$ so that $\varphi(t+1, \mathbf{b}) \leq \mathbb{E}_{t+1}[f_{jl}Q_{jl}(t+1) + \varphi(t, \mathbf{b} - A_j Y_{jl}^{t+1})]$, implying that posting price f_{jl} 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 $\mathcal{Q}(t, B^t)$, and bound the information loss $\sum_t \mathbb{P}[\mathcal{Q}(t, B^t)]$.

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 $\mathbb{P}[\mathcal{Q}(t, B^t)] \leq 1/t^2$ 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 Appendix C. 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 $n = 1$, i.e., selling multiple copies of an item to homogeneous customers; since there is only one type, we drop the index j . Recall $f_1 > \dots > f_m$ and $q_1 < \dots < q_m$. It is easy to check that the solution of $P[b, \mathbf{q}, t]$ is as follows: (i) If $b \leq tq_1$, then $x = (b/q_1, 0, \dots, 0)$; (ii) If $b > tq_m$, then $x = (0, \dots, 0, t)$; (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; 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 , $\mathbb{P}\left[\max_l |Q_{jl}(t) - \mathbb{E}[Q_{jl}(t)]| > \sqrt{\frac{\log(t)}{t}}\right] \leq \frac{c_j}{t^2}$.

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

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

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

$$\mathbb{P}\left[\sup_l |Q_{jl}(t) - \bar{F}(f_{jl})| > \sqrt{\frac{\log(t)}{t}}\right] \leq 2(1 - p_j + p_j e^{-\theta})^t \quad \text{where } \theta = 2\log(t)/t.$$

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

Stability of Left-Hand Side Perturbations. As stated in Algorithm 3, ONLINE takes actions based on $P[\mathbf{b}, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]]$, while 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 Δ 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 $\Delta\mathbf{z} = Z^t - \mathbb{E}[Z^t]$.

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

$$P_S[V_t, \mathbf{q} + \Delta\mathbf{q}, \mathbf{z} + \Delta\mathbf{z}] := \max \left\{ x_{j'l'} : \sum_{j,l} f_{jl}(q_{jl} + \Delta q_{jl}) x_{jl} \geq V_t, \mathbf{x} \text{ feasible for } P[\mathbf{b}, \mathbf{q} + \Delta\mathbf{q}, \mathbf{z} + \Delta\mathbf{z}] \right\}.$$

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

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) Δ , 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 below; the proof is provided in Appendix C.

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 $t \geq c$, with probability $1 - c/t^2$, $P_S[V_t, Q(t), Z^t] - P_S[v_t, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]] \geq -c\sqrt{t \log(t)}$.

Lemma 7 leads to the bound in Proposition 2. Indeed, since the LP in Eq. (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, $P_S[v_t, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]] \geq tp_j/(m+1)$. We know that posting f_{jv} is satisfying whenever $P_S[V_t, Q(t), Z^t] \geq 1$ (see Lemma 6), hence posting the maximum entry is satisfying provided that $tp_j/(m+1) - c\sqrt{t \log(t)} \geq 1$, 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 multi-dimensional 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 vs. RABBI gap.

For the small system, the k -th 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 large system, the parameters were generated randomly and are reported in Appendix F, the k -th 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 re-solving an LP in each period). In Fig. 5

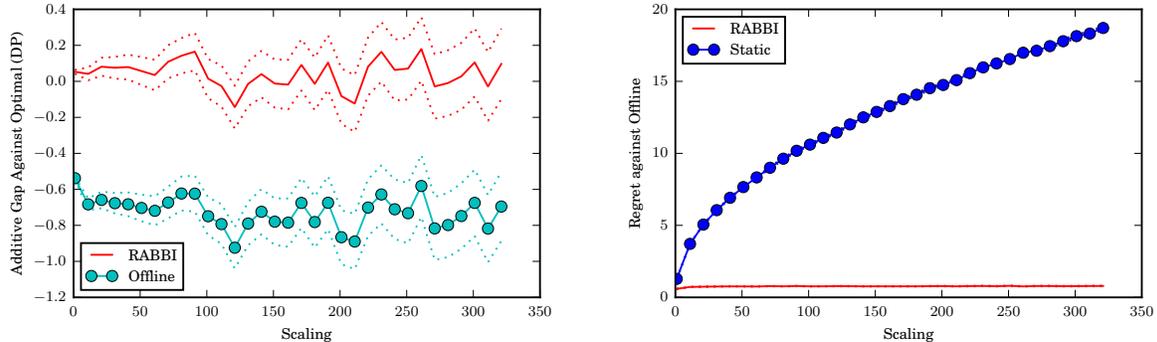


Figure 5 Regret in the ‘small system’ ($n = 1$ and $d = 1$), with horizon $T = 20k$ and initial budget $B = 6k$, under scaling $k = 1, 10, 20, \dots, 340$. Dotted lines represent 90% CI. (LEFT) additive gaps against the optimal policy, i.e., $V^{DP} - V^{\text{RABBI}}$ and $V^{DP} - V^{\text{OFFLINE}}$ (RIGHT) Regret of two policies against OFFLINE.

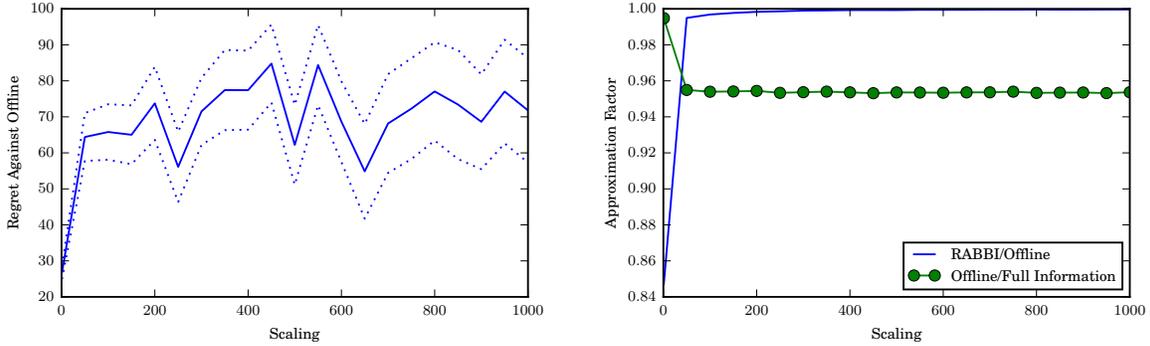


Figure 6 Performance in the ‘large system’ ($n = 20$ and $d = 25$) with horizon $T = 100k$ and initial budgets $B_i = 10k$ for $i \in [25]$, under scaling $k = 1, 2, \dots, 1000$. (LEFT) Regret against OFFLINE, i.e., $V^{\text{OFFLINE}} - V^{\text{RABBI}}$; dotted lines represent 90% CI. (RIGHT) Approximation factors $V^{\text{RABBI}}/V^{\text{OFFLINE}}$ of RABBI compared to 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.

(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 Fig. 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 $f \cdot D(f) = f \cdot T \cdot \bar{F}(f)$ subject to the constraint $D(f) \leq B$. The solution is the better of two prices: (i) the market clearing price, i.e., that satisfies $D(f) = B$ or (ii) the monopoly price which maximizes $fD(f)$. We note though that when a continuum of prices are allowed, (Jasin 2014) propose an algorithm (that, like RABBI, is based on resolving an optimization problem in each period) which achieves a regret that is logarithmic in k under certain non-degeneracy 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 competitive ratio of *any non-anticipatory 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, make a choice 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 $(f_{1\alpha}, \dots, f_{n\alpha})$, 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 α , 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 no-purchase option). These probabilities might be derived, for example, from a standard family such as the multinomial-logit model, nested logit model, etc.; 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 $B \in \mathbb{N}^d$.*

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

$$(P[\mathbf{b}, \mathbf{q}, \mathbf{z}]) \quad \text{maximize:} \quad \sum_{\alpha \in \mathcal{A}} x_{\alpha} \sum_{j \in [n]} f_{j\alpha} q_{j\alpha} \quad (13)$$

$$\begin{aligned} \text{subject to: } \quad & \sum_{\alpha \in \mathcal{A}} \sum_{j \in [n]} a_{ij} q_{j\alpha} x_{\alpha} \leq b_i \quad \forall i \in [d] \\ & \sum_{\alpha \in \mathcal{A}} x_{\alpha} = t \\ & \mathbf{x} \geq 0 \end{aligned}$$

Here $q_{j\alpha}$ stands for the fraction of customers that would buy product j if presented with the price assortment α . 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 α , 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 Eq. (13) with $q_{j\alpha} = Q_{j\alpha}(t)$, where :

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

With the (re)defined key ingredients—namely the LP in Eq. (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 α is posted at time t , the random collected reward is $\sum_j Y_{j\alpha}^t f_{j\alpha}$ and the random inventory at $t-1$ is $b - \sum_j A_j Y_{j\alpha}^t$. In turn, the Bellman loss in Eq. (12) takes on the form

$$L_B(t+1, \mathbf{b}, \alpha) = P[\mathbf{b} - \sum_j A_j Q_{j\alpha}(t+1), Q(t+1), t] - \mathbb{E}_{t+1}[P[\mathbf{b} - \sum_j A_j Y_{j\alpha}^{t+1}, Q(t), t]] \quad (14)$$

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

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 7 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 Fig. 7.

	Products	High	High-discount	Low	Low-discount	Only R2	Only R1	Stock-out
$f_{j\alpha}$	{R1}	5	5	3	3	∞	5	Out
	{R2}	5	5	3	3	5	∞	∞
	{R1,R2}	10	9	6	5	∞	∞	∞
$p_j(\alpha)$	{R1}	0.2	0.2	0.3	0.3	0	0.2	0
	{R2}	0.2	0.2	0.3	0.3	0.2	0	0
	{R1,R2}	0.1	0.15	0.2	0.25	0	0	0

Table 1 Example with seven assortments: 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.

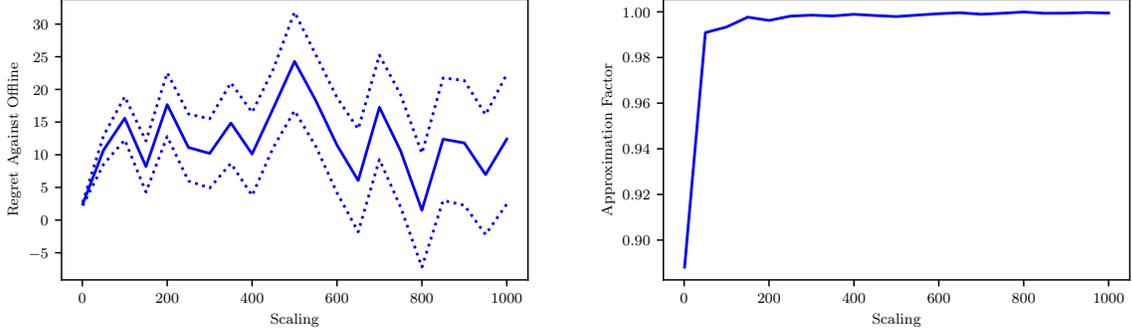


Figure 7 Performance of pricing-RABBI with customer choice (see Table 1). We set the horizon as $T = 10k$ and the inventory as $(R1, R2) = (3k, 2k)$, and vary scaling parameter $k = 1, \dots, 1000$. (LEFT) Regret against OFFLINE, with 90% confidence intervals. (RIGHT) Approximation ratio of RABBI against DP.

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 of 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 $t \in \{T, T-1, \dots, 1\}$ 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 $\mathcal{S} = \mathbb{R}_{\geq 0} \times \{\emptyset, \mathbf{a}, \mathbf{r}\}$, 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 $\diamond \in \{\mathbf{a}, \mathbf{r}\}$, reducing the capacity if $\diamond = \mathbf{a}$. At the second stage, the state is of the form $s = (b, \diamond)$ with $\diamond \in \{\mathbf{a}, \mathbf{r}\}$, and we collect the reward only if $\diamond = \mathbf{a}$. Formally, the rewards are $\mathcal{R}((b, \mathbf{a}), \xi^{t-1/2}, \emptyset) = \xi^{t-1/2}$ and $\mathcal{R}((b, \mathbf{r}), \xi^{t-1/2}, \emptyset) = 0$.

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

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

$$(P[b, \mathbf{y}, \mathbf{z}]) \quad \text{maximize:} \quad \sum_j y_j x_j \quad (15)$$

$$\begin{aligned}
\text{subject to: } \quad & \sum_j w_j x_{j\mathbf{a}} \leq b \\
& x_{j\mathbf{a}} + x_{j\mathbf{r}} = z_j \quad \forall j \in [n] \\
& \mathbf{x} \geq 0
\end{aligned}$$

Note that if the average rewards \mathbf{r} were known, then setting $\mathbf{y} = \mathbf{r}$ we get the LP relaxation of Eq. (1) for the baseline `OnlineKnapsack`. Moreover, for any \mathbf{r} , the optimal LP solution sorts types by their ‘‘bang for the buck’’ ratios r_j/w_j , and accepts them greedily. In particular, the solution only requires knowing the ranking induced by \mathbf{r} .

Offline Benchmark and Relaxed Value Function: In this setting, we define `OFFLINE` as the controller that knows the number of arrivals Z_j^T for each j , and also, *knows the ranking of the types* (i.e., knows $r_j/w_j \forall j \in [n]$).

Formally, `OFFLINE` is defined via the filtration $\mathcal{G}_t = \sigma(\{\xi^t : t \in [T]\} \cup \{\xi^\tau : \tau \geq t\})$. 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

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

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 an `ONLINE` control problem with distribution learning in this setting. \blacksquare

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.

6.2. Regret Analysis for Distribution-Agnostic Online Knapsack

As in the earlier sections, we first demonstrate that φ satisfies the Bellman inequalities

LEMMA 8. *The relaxation φ defined in (16) satisfies the Bellman Inequalities with exclusion sets*

$$\mathcal{B}(t, b) = \{\omega \in \Omega : \exists X \text{ solving } (P[b, \mathbf{r}, Z^t]) \text{ s.t. } X_{\xi^t, \mathbf{a}} \geq 1 \text{ or } X_{\xi^t, \mathbf{r}} \geq 1\}.$$

PROOF. The initial ordering in Definition 2 follows from an argument identical to that of Lemma 2. The monotonicity property follows from Proposition 3. \square

Algorithm 4 Learning RABBI**Input:** Access to solutions of $(P[b, \mathbf{y}, \mathbf{z}])$ **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** $t \in \{T, T-1, \dots, 1\}$ **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 $\hat{U}^t \in \operatorname{argmax}_{u=\mathbf{a}, \mathbf{r}} \{X_{ju}^t\}$
- 5: If $\hat{U}^t = \mathbf{a}$, 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_j .

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 $r_\varphi(\operatorname{Regret}_1 + \operatorname{Regret}_2)$, where Regret_1 and Regret_2 are two specific sources of error. When the estimators R^t of \mathbf{r} are accurate enough, the error is Regret_1 and is attributed to the incorrect “guess” of a satisfying action, i.e., Regret_1 is an algorithmic regret. The second term, Regret_2 , is the error that arises from insufficient accuracy of R^t , i.e., 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

$$\operatorname{Regret}_1 \leq 2 \sum_j \frac{(w_{\max}/w_j)^2}{p_j} \quad \text{and} \quad \operatorname{Regret}_2 \leq 16 \sum_j \frac{1}{p_j (w_j \delta)^2}.$$

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

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) $\mathbb{P}[|Z^t - \mathbb{E}[Z^t]| \geq c\mathbb{E}[Z^t]]$ 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. \blacksquare

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 $\delta > 0$: $|\bar{r}_j - \bar{r}_{j'}| \geq \delta$ for all $j \neq j'$, where $\bar{r}_j = \mathbb{E}[R_j]/w_j$.

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 is 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), ε -Greedy or simply to set apart some time for exploration (see Corollary 1 below).

Recall that $\sigma : [n] \rightarrow [n]$ is the ordering of $[n]$ w.r.t. the ratios $\bar{r}_j = r_j/w_j$ and $\hat{\sigma}^t : [n] \rightarrow [n]$ 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 $\{\mathbf{explore}, \mathbf{exploit}\}$. If the action is $\mathbf{explore}$, 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, i.e., 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** $t = T, \dots, 1$ **do**
 - 3: Observe input ξ^t and let $U \leftarrow \text{BANDITS}(T, t, S^t, \mathcal{F}_t)$.
 - 4: If $U = \mathbf{explore}$, accept the arrival
 - 5: If $U = \mathbf{exploit}$, 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_1 be the regret of Algorithm 4, as given in Theorem 4. Define the indicators $\mathbf{explore}_t, \mathbf{exploit}_t$ which denote the output of BANDITS at time t . The regret of Algorithm 5 is at most $r_\varphi M$, where*

$$M = \text{Regret}_1 + \mathbb{E} \left[\sum_t \mathbf{explore}_t \right] + \mathbb{E} \left[\sum_t \mathbb{P}[\sigma \neq \hat{\sigma}^t] \mathbf{exploit}_t \right].$$

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 $\mathbf{explore}$ (which, in particular, implies accepting the item), the decision might be the wrong one (i.e., different than 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 $\mathbf{exploit}$ (in which case we call Algorithm 4), but the (learned) ranking at time t , $\hat{\sigma}^t$, is different than σ^t . This is the last term in M . Finally, even if the learned ranking is correct, $\mathbf{exploit}$ 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 $\frac{8}{(w_j\delta)^2} \log T$ 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 re-solving algorithms. The framework is based on using a carefully chosen offline benchmark, that 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 φ satisfying the Bellman Inequalities requires some ad-hoc creativity but, as our example illustrate, is often rather intuitive. In Appendix A we provide sufficient conditions, applicable to cases where φ 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 φ , 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.

Acknowledgments

SB and AV gratefully acknowledge support from the ARL under grant W911NF-17-1-0094, and the NSF under grants CNS-1955997, DMS-1839346 and ECCS-1847393; IG’s work was supported by the DoD under grant W911NF-20-C-0008.

References

- Agrawal S, Devanur N (2016) Linear contextual bandits with knapsacks. *Advances in Neural Information Processing Systems*, 3450–3458.
- Alaei S (2014) Bayesian combinatorial auctions: Expanding single buyer mechanisms to many buyers. *SIAM Journal on Computing* 43(2):930–972.
- Arlotto A, Gurvich I (2019) Uniformly bounded regret in the multisecretary problem. *Stochastic Systems* 9(3):231–260.

-
- Babai M, Dughmi S, Kleinberg R, Slivkins A (2015) Dynamic pricing with limited supply. *ACM Transactions on Economics and Computation (TEAC)* 3(1):4.
- Badanidiyuru A, Kleinberg R, Slivkins A (2018) Bandits with knapsacks. *Journal of the ACM (JACM)* 65(3):13.
- Badanidiyuru A, Langford J, Slivkins A (2014) Resourceful contextual bandits. *Conference on Learning Theory*, 1109–1134.
- Balseiro SR, Brown DB (2019) Approximations to stochastic dynamic programs via information relaxation duality. *Operations Research* 67(2):577–597.
- Banerjee S, Freund D (2020) Uniform loss algorithms for online stochastic decision-making with applications to bin packing. *ACM SIGMETRICS*.
- Boucheron S, Lugosi G, Massart P (2013) *Concentration inequalities: A nonasymptotic theory of independence* (Oxford university press).
- Brown D, Smith J, Sun P (2010) Information Relaxations and Duality in Stochastic Dynamic Programs. *Operations Research* 58(4):785–801.
- Bubeck S, Cesa-Bianchi N, et al. (2012) Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning* .
- Bubeck S, Perchet V, Rigollet P (2013) Bounded regret in stochastic multi-armed bandits. *Conference on Learning Theory*, 122–134.
- Buchbinder N, Jain K, Singh M (2014) Secretary problems via linear programming. *Mathematics of Operations Research* 39(1):190–206.
- Bumpensanti P, Wang H (2020) A re-solving heuristic for dynamic resource allocation with uniformly bounded revenue loss. *Management Science* .
- Chen Q, Jasin S, Duenyas I (2019) Nonparametric self-adjusting control for joint learning and optimization of multiproduct pricing with finite resource capacity. *Mathematics of Operations Research* .
- Chugg B, Maehara T (2019) Submodular stochastic probing with prices. *International Conference on Control, Decision and Information Technologies (CoDIT)*, 60–66 (IEEE).
- Correa J, Foncea P, Hoeksma R, Oosterwijk T, Vredeveld T (2017) Posted price mechanisms for a random stream of customers. *ACM EC*, 169–186.
- Dubhashi D, Panconesi A (2009) *Concentration of measure for the analysis of randomized algorithms* (Cambridge University Press).
- Düetting P, Feldman M, Kesselheim T, Lucier B (2017) Prophet inequalities made easy: Stochastic optimization by pricing non-stochastic inputs. *IEEE FOCS*.
- Gallego G, Van Ryzin G (1997) A multiproduct dynamic pricing problem and its applications to network yield management. *Operations research* 45(1):24–41.

- Gupta A, Nagarajan V (2013) A stochastic probing problem with applications. *International Conference on Integer Programming and Combinatorial Optimization*.
- Gupta A, Nagarajan V, Singla S (2016) Algorithms and adaptivity gaps for stochastic probing. *ACM-SIAM SODA* .
- Jasin S (2014) Reoptimization and self-adjusting price control for network revenue management. *Operations Research* 62(5):1168–1178.
- Jasin S, Kumar S (2012) A re-solving heuristic with bounded revenue loss for network revenue management with customer choice. *Mathematics of Operations Research* 37(2):313–345.
- Kleinberg R, Weinberg SM (2012) Matroid prophet inequalities. *ACM STOC*.
- Mangasarian O, Shiao T (1987) Lipschitz Continuity of Solutions of Linear Inequalities, Programs and Complementarity Problems. *SIAM Journal on Control and Optimization* 25(3):583–595.
- Massart P (1990) The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability* 1269–1283.
- Shank N, Yang H (2013) Coupon Collector Problem for Non-Uniform Coupons and Random Quotas. *The Electronic Journal of Combinatorics* 20(2):33.
- Singla S (2018) The price of information in combinatorial optimization. *ACM-SIAM SODA*.
- Talluri KT, Van Ryzin GJ (2006) *The theory and practice of revenue management*, volume 68 (Springer Science & Business Media).
- Vera A, Banerjee S (2020) The Bayesian Prophet: A Low-Regret Framework for Online Decision Making. *Management Science* .
- Weitzman ML (1979) Optimal search for the best alternative. *Econometrica: Journal of the Econometric Society* 641–654.
- Wu H, Srikant R, Liu X, Jiang C (2015) Algorithms with logarithmic or sublinear regret for constrained contextual bandits. *Advances in Neural Information Processing Systems*, 433–441.

Appendix A: A Sufficient Condition for Bellman Inequalities

In this section, we construct φ based on a general optimization program and provide a sufficient condition to guarantee monotonicity. This serves to underscore some of the key elements in a problem's structure that allows one to construct low regret online policies. This guideline does not apply to all the examples we study here: in particular, it applies to the baseline and learning variants, but not to probing or pricing.

We study a particular case of canonical filtrations (see Definition 1), where the random variables G_θ that we reveal are the inputs ξ^θ for some fixed times Θ (see Fig. 2 for an illustration).

Recall that we reveal some inputs to OFFLINE, but not necessarily all of them; we call *concealed inputs* those not revealed to OFFLINE. Informally speaking, we will show that φ satisfies the Bellman Inequalities if (i) OFFLINE's relaxed value φ can be computed with a linear program and (ii) the concealed inputs are in the objective function only (not in the constraints). Requirements (i) and (ii) are appealing because they are verifiable directly from the problem structure without any computation.

Recall that, with t periods to go, OFFLINE knows the randomness $\{\xi^T, \dots, \xi^t, \xi_\Theta\}$, where we denote $\xi_\Theta = (\xi^\theta : \theta \in \Theta)$. In other words, we reveal $\{\xi^T, \dots, \xi^t, \xi_\Theta\}$, while the inputs $\{\xi^l : l < t, l \notin \Theta\}$ are concealed.

Suppose the relaxation is an LP with decision variables \mathbf{x} (see Eq. (5)):

$$\varphi(t, s | \mathcal{G}_t) = \max_{\mathbf{x} \in \mathbb{R}^{\Xi \times [T] \times \mathcal{U}}} \{\mathbb{E}[h(\mathbf{x}; \xi^1, \dots, \xi^T) | \mathcal{G}_t] : g(\mathbf{x}; s, \xi^T, \dots, \xi^t, \xi_\Theta) \geq \mathbf{0}\}, \quad (17)$$

where \mathcal{U}, Ξ are the control and input spaces. For input ξ , control u , and time t , we interpret $x_{\xi, t, u}$ as a variable indicating if OFFLINE uses u at time t when presented input ξ .

PROPOSITION 3. *Let h, g be linear functions and let φ be given by (17). Assume further that the following holds for all s, t, u*

(i) *h captures rewards: $\mathbb{E}[h(\mathbf{e}_{\xi^t, t, u}; \xi^1, \dots, \xi^T) | \mathcal{G}_t] \leq \mathcal{R}(s, \xi^t, u)$ for actions u that are feasible in state s .*

(ii) *g captures transitions: $g(\mathbf{e}_{\xi^t, t, u}; s, \xi^T, \dots, \xi^t, \xi_\Theta) \leq g(\mathbf{0}; \mathcal{T}(s, \xi^t, u), \xi^T, \dots, \xi^{t-1}, \xi_\Theta)$.*

Then, φ satisfies monotonicity with exclusion sets

$$\mathcal{B}(t, s) = \{\omega \in \Omega : \exists X[\omega] \text{ solving } \varphi(t, s | \mathcal{G}_t) \text{ s.t. } X_{\xi^t, t, u} \geq 1 \text{ for some } u \in \mathcal{U}\}.$$

It is natural to say that h captures the reward if the incremental effect of taking the action u given input ξ^t is equal to the immediate reward $\mathbb{E}[h(\mathbf{e}_{\xi^t, t, u}; \xi^1, \dots, \xi^T) | \mathcal{G}_t] = \mathcal{R}(s, \xi^t, u)$. It is similarly natural to say that g captures transitions if it is *stable under the one-step transition*, namely, that $g(\mathbf{e}_{\xi^t, t, u}; s, \xi^T, \dots, \xi^t, \xi_\Theta) = g(\mathbf{0}; \mathcal{T}(s, \xi^t, u), \xi^T, \dots, \xi^{t-1}, \xi_\Theta)$; in other words, this means that taking

the action u at time t , has the same effect as taking no action at the state $\mathcal{T}(s, \xi^t, u)$. This should hold in any reasonable resource consumption problem, e.g., consuming 1 with B units of budget remaining is the same as not consuming anything with $B - 1$ units. In the result below we make the weaker assumption that these relationships hold as inequalities.

The baseline and learning variants are useful illustrations of Proposition 3.

EXAMPLE 3 (BASELINE). Let \mathcal{G} be the full information filtration ($\Theta = [T]$). In Section 2.2 we introduced a linear relaxation for OFFLINE. We start by writing a relaxation in the form of Proposition 3 and show how it subsequently simplifies to the final form in Section 2.2.

Recall that \mathbf{a}, \mathbf{r} denote the actions accept and reject. A natural “expanded” linear program is

$$\max \left\{ \sum_j \sum_{l=1}^t x_{j,l,\mathbf{a}} r_j : \sum_{j,l} w_j x_{j,l,\mathbf{a}} \leq s, 0 \leq x_{j,l,\mathbf{a}} \leq \mathbb{1}_{\{\xi^l=j\}} \right\}.$$

Defining the auxiliary variables $x_j := \sum_l x_{j,l,\mathbf{a}}$, this is equivalent to $\varphi(t, s | \mathcal{G}) = \max\{\mathbf{r}'\mathbf{x} : \mathbf{w}'\mathbf{x} \leq s, 0 \leq \mathbf{x} \leq \mathbf{Z}^t\}$, where, recall $Z_j^t = \sum_{l=1}^t \mathbb{1}_{\{\xi^l=j\}}$ counts the number of type- j arrivals in the last t periods.

This φ also has the form of Proposition 3, with the functions h and g given by (note that the action \mathbf{r} has zero objective coefficient)

$$h(\mathbf{x}; \xi^1, \dots, \xi^t) := \sum_j x_{j\mathbf{a}} r_j \quad \text{and} \quad g(\mathbf{x}; s, \xi^T, \dots, \xi^1) := \left(\begin{array}{c} s - \sum_j x_{j\mathbf{a}} \\ Z^t - \mathbf{x} \end{array} \right).$$

Conditions (i) and (ii) can be easily verified now. The objective h is a linear function of the decision vector \mathbf{x} and the constraint function g aggregates ξ into the sums Z^t . \blacksquare

In the learning setting, OFFLINE is presented with a public type j and must decide whether to accept or reject before seeing the private type, which is a reward R_j drawn from an unknown distribution.

EXAMPLE 4 (LEARNING). Let us model the problem with $2T$ time periods, where at even times the public type is revealed and at odd times the private (reward). In this model, the input ξ^t is an index $j \in [n]$ at even times and it is a reward $R \in \mathbb{R}$ at odd times. Also let us model the random rewards by drawing i.i.d. copies $\{R_{jt}\}_t$ of R_j .

Let us endow OFFLINE with the information of all even times, i.e., OFFLINE knows all the future arriving public types. Specifically, we set $\Theta = \{t \in [T] : t \text{ is even}\}$ (see Fig. 2 for a representation of \mathcal{G}). The realizations $\{R_{jt}\}_{j,t}$, drawn at times $t \notin \Theta$, are concealed. The expanded linear program is

$$\max \left\{ \sum_j \sum_{l=1}^t x_{j,l,\mathbf{a}} \mathbb{E}[R_j] : \sum_{j,l} w_j x_{j,l,\mathbf{a}} \leq s, 0 \leq x_{j,l,\mathbf{a}} \leq \mathbb{1}_{\{\xi^l=j\}} \right\}.$$

As before, we can simplify this LP by aggregating variables, see Section 6 for the details. Here we prefer to study the expanded LP because it exemplifies the conditions in Proposition 3.

The objective function is $h(\mathbf{x}; \xi^1, \dots, \xi^t) = \sum_{j,l} x_{j,l,a} R_{j,l}$. When we take expectations $\mathbb{E}[\cdot | \mathcal{G}_t]$ we arrive at the expression $\sum_{j,l} x_{j,l,a} \mathbb{E}[R_{j,l}]$. The constraint function g is given by the feasibility region of the LP. Conditions (i) and (ii) of Proposition 3 hold with equality. ■

PROOF OF PROPOSITION 3. Let $u \in \mathcal{U}$ be such that $X_{\xi^t, t, u} \geq 1$. Denote $\theta_t := \{l \in [T] : l \geq t\} \cup \Theta$, so all the inputs $(\xi^l : l \in \theta_t)$ are revealed at time t (the rest are concealed). By Lemma 1,

$$\varphi(t, s | \mathcal{G}_t) = \mathbb{E}[h(\mathbf{e}_{\xi^t, t, u}; \xi^1, \dots, \xi^T) | \mathcal{G}_t] + \max_{\mathbf{x}} \{\mathbb{E}[h(\mathbf{x}; \xi^1, \dots, \xi^T) | \mathcal{G}_t] : g(\mathbf{x} + \mathbf{e}_{\xi^t, t, u}; s, (\xi^l : l \in \theta_t)) \geq \mathbf{0}\}.$$

Using (i) and (ii) yields

$$\varphi(t, s | \mathcal{G}_t) \leq \mathcal{R}(s, \xi^t, u) + \max_{\mathbf{x}} \{\mathbb{E}[h(\mathbf{x}; \xi^1, \dots, \xi^T) | \mathcal{G}_t] : g(\mathbf{x}; \mathcal{T}(s, \xi^t, u), (\xi^l : l \in \theta_{t-1})) \geq \mathbf{0}\}. \quad (18)$$

Since \mathcal{G}_t is coarser than \mathcal{G}_{t-1} , we know that $\mathbb{E}[\mathbb{E}[\cdot | \mathcal{G}_{t-1}] | \mathcal{G}_t] = \mathbb{E}[\cdot | \mathcal{G}_t]$. Using Eq. (18) and applying Jensen's Inequality (recall that the maximum of linear functions is a convex function) we obtain

$$\varphi(t, s | \mathcal{G}_t) \leq \mathcal{R}(s, \xi^t, u) + \mathbb{E} \left[\max_{\mathbf{x}} \{\mathbb{E}[h(\mathbf{x}; \xi^1, \dots, \xi^T) | \mathcal{G}_{t-1}] : g(\mathbf{x}; \mathcal{T}(s, \xi^t, u), (\xi^l : l \in \theta_{t-1})) \geq \mathbf{0}\} \middle| \mathcal{G}_t \right].$$

This corresponds to the required inequality in Definition 2. □

The sufficient conditions in Proposition 3 are not necessary; they are not satisfied in the probing setting (Section 4) or in the pricing setting (Section 5). Nevertheless, we are still able to show monotonicity and draw the desired regret bounds.

Appendix B: Additional Details from Section 4 (Online Probing)

We first state and prove an auxiliary lemma which we need for our proofs.

LEMMA 9. Consider the standard-form LP $(P[\mathbf{d}]) : \max\{\mathbf{r}'\mathbf{x} : M\mathbf{x} = \mathbf{d}, \mathbf{x} \geq 0\}$, where $M \in \mathbb{R}^{m \times n}$ is an arbitrary constraint matrix and $\mathbf{d} \in \mathbb{R}^m$. The function $\mathbf{d} \mapsto P[\mathbf{d}]$ is concave and therefore, if X is a random right-hand side, then $\mathbb{E}[P[X]] \leq P[\mathbb{E}[X]]$.

PROOF. The dual problem is $(D[\mathbf{d}]) : \min\{\mathbf{d}'\mathbf{y} : M'\mathbf{y} \geq \mathbf{r}\}$. The function $\mathbf{d} \mapsto D[\mathbf{d}]$ is a minimum of linear functions, therefore concave. □

B.1. Bellman Inequalities and Loss

We first establish the initial ordering property.

PROOF OF LEMMA 2. Consider a policy for OFFLINE determining when to probe, accept or reject. Recall such a policy is a mapping $\pi : [T] \times \mathcal{S} \rightarrow \mathcal{U}$ s.t. $\pi(t, s)$ is \mathcal{G}_t -measurable for all t, s .

The policy, once fixed, induces a random trajectory determined by the realization of the probed rewards. Denote the random number of times where a type j was probed as X_{jp} , accepted (rejected) without probing as X_{ja} (X_{jr}), and accepted (rejected) after probe outcome is (j, k) as X_{jka} (X_{jkr}). Then, we can write $\mathbb{E}[V(T, (b_h, b_p)|\mathcal{G}_T)] = \mathbb{E}[\sum_j \bar{r}_j X_{ja} + \sum_{j,k} r_{jk} X_{jka} | \mathcal{G}_T]$, where we use the fact that, conditional on accepting without probing, the expected reward is \bar{r}_j . Thus we have

$$\mathbb{E}[V(t, (b_h, b_p)|\mathcal{G}_T)] = \sum_j \bar{r}_j \mathbb{E}[X_{ja} | \mathcal{G}_T] + \sum_{j,k} r_{jk} \mathbb{E}[X_{jka} | \mathcal{G}_T].$$

We now claim that $\mathbb{E}[\mathbf{X}]$ yields a feasible solution to $(P[T, (b_h, b_p), Z])$. Indeed, with the exception of the constraint $x_{jka} + x_{jkr} = q_{jk} x_{jp}$, the random variables satisfy a.s. all the constraints of $(P[T, (b_h, b_p), Z])$. Furthermore, since OFFLINE's policy is adapted to \mathcal{G} , we obtain $\mathbb{E}[X_{jka} + X_{jkr} | X_{jp}, \mathcal{G}_T] = q_{jk} X_{jp}$, thus the expected values satisfy the desired constraint. To summarize, $V(T, (b_h, b_p)|\mathcal{G}_T)$ equals the value of the feasible solution given by the expectations. \square

Next we establish the monotonicity condition in Definition 2.

PROOF OF LEMMA 3. Observe that the monotonicity condition in Definition 2 translates to the following condition in the online probing setting.

$$\varphi(t, (b_h, b_p, \emptyset) | \mathcal{G}_t) \leq \max_{\diamond \in \{\mathbf{a}, \mathbf{p}, \mathbf{r}\}} \{ \mathbb{E}_{\xi^{t-1/2}} [\varphi(t-1/2, (s_\diamond, \diamond) | \mathcal{G}_{t-1/2}) | \mathcal{G}_t] \} \quad \forall \omega \notin \mathcal{B}(t, s).$$

where the state $s_\diamond = (b_h - 1, b_p)$ if $\diamond = \mathbf{a}$, $s_\diamond = (b_h, b_p - 1)$ if $\diamond = \mathbf{p}$ and $s_\diamond = (b_h, b_p)$ if $\diamond = \mathbf{r}$.

First, given $\xi^t = i$, we have from Eq. (8) that $\mathbb{E}_{\xi^{t-1/2}} [\varphi(t-1/2, (s_\diamond, \diamond) | \mathcal{G}_{t-1/2}) | \mathcal{G}_t] = P[(b_h, b_p), Z^{t-1}]$ if $\diamond = \mathbf{r}$, and $r_{\xi^{t-1/2}} + P[(b_h - 1, b_p), Z^{t-1}]$ if $\diamond = \mathbf{a}$. Now for cases (1) and (2), the claim in the lemma follows directly by invoking Lemma 1.

For case (3) we need to introduce some notation. Let $\mathbf{q}_j \in \mathbb{R}^{n \times m}$ be a vector with value q_{jk} in components (j, k) , $k \in [m]$, and zero otherwise (i.e. in components (j', k) with $j' \neq j$). Similarly, let $\mathbf{1}_{(j,k)} \in \mathbb{R}^{n \times m}$ have value 1 in the single component (j, k) and zero otherwise. We also rewrite the LP in Eq. (7) with an extra 'budget vector' \mathbf{w} such that $P[(b_h, b_p), \mathbf{z}] = \bar{P}[(b_h, b_p), \mathbf{z}, \mathbf{0}]$.

$$\begin{aligned} (\bar{P}[(b_h, b_p), \mathbf{z}, \mathbf{w}]) \quad & \text{maximize:} \quad \sum_{j,k} r_{jk} x_{jka} + \sum_j \bar{r}_j x_{ja} \\ & \text{subject to:} \quad \sum_{j,k} x_{jka} + \sum_j x_{ja} \leq b_h \\ & \quad \quad \quad \sum_j x_{jp} \leq b_p \\ & \quad \quad \quad x_{ja} + x_{jp} + x_{jr} = z_j \quad \forall j \in [n] \\ & \quad \quad \quad x_{jka} + x_{jkr} - q_{jk} x_{jp} = w_{jk} \quad \forall j \in [n], k \in [m] \\ & \quad \quad \quad \mathbf{x} \geq \mathbf{0} \end{aligned}$$

Now if $\bar{X}_{ip} \geq 1$ and $\xi^{t-1/2} = (i, k)$ is such that either $\bar{X}_{ika} \geq 1$ or $\bar{X}_{ikr} \geq 1$, then by Lemma 1, we have the following decomposition (depending on the random $\xi^{t-1/2}$)

$$\bar{P}[(b_h, b_p), Z^t, \mathbf{0}] = r_{\xi^{t-1/2}} \mathbf{1}_{\{\bar{X}_{ika} \geq 1\}} + \bar{P}[(b_h - \mathbf{1}_{\{\bar{X}_{ika} \geq 1\}}, b_p - 1), Z^{t-1}, \mathbf{q}_{\xi^t} - \mathbf{1}_{\xi^{t-1/2}}], \quad \forall \omega \notin \mathcal{B}(t, b_h, b_p)$$

where the vectors $\mathbf{q}, \mathbf{1}$ are evaluated in random components; since by assumption $\bar{X}_{ip} \geq 1$ under the optimal solution, the optimal value in the optimization problem is the same as the reward obtained “now” ($r_{\xi^{t-1/2}}$) and the residual value after discounting b_p by one. Taking expectations $\mathbb{E}[\cdot | \mathcal{G}_t]$ and using Lemma 9 we have

$$\begin{aligned} P[(b_h, b_p), Z^t] &= \mathbb{E}[r_{\xi^{t-1/2}} \mathbf{1}_{\{\bar{X}_{ika} \geq 1\}} + \bar{P}[(b_h - \mathbf{1}_{\{\bar{X}_{ika} \geq 1\}}, b_p - 1), Z^{t-1}, \mathbf{q}_{\xi^t} - \mathbf{1}_{\xi^{t-1/2}}] | \mathcal{G}_t] \\ &\leq \mathbb{E}[r_{\xi^{t-1/2}} \mathbf{1}_{\{\bar{X}_{ika} \geq 1\}} + \bar{P}[(b_h - \mathbf{1}_{\{\bar{X}_{ika} \geq 1\}}, b_p - 1), Z^{t-1}, \mathbf{0}] | \mathcal{G}_t] \\ &\leq \mathbb{E}[\max\{r_{\xi^{t-1/2}} + P[(b_h - 1, b_p - 1), Z^{t-1}], P[(b_h, b_p - 1), Z^{t-1}] | \mathcal{G}_t]. \end{aligned}$$

The last inequality, following from substituting $\mathbf{1}_{\{\bar{X}_{ika} \geq 1\}} \in \{0, 1\}$, gives the desired result. \square

Appendix C: Additional Details from Section 5 (Dynamic Pricing)

C.1. Proof of Lemma 4

Throughout this subsection, we fix some indexes j', l' . To complete the proof of the proposition, it remains to establish that, whenever $X_{j'l'} \geq 1$, then $\mathbb{E}[L_B(t+1, \mathbf{b}, j', l') | \mathcal{G}_t] \leq 0$, where

$$L_B(t+1, \mathbf{b}, j', l') = P[\mathbf{b} - A_{j'} Q_{j'l'}(t+1), Q(t+1), Z^t] - \mathbb{E}_{t+1}[P[\mathbf{b} - A_{j'} Y_{j'l'}, Q(t), Z^t]].$$

The Correction LP. Let us fix $(t, \mathbf{b}, \mathbf{q}, \mathbf{z})$ and denote $\bar{\mathbf{x}}$ the solution of $P[\mathbf{b} - A_{j'} q_{j'l'}, \mathbf{q}, \mathbf{z}]$. To bound the loss, we must bound the right-hand side of (12), which captures the perturbation of budgets from $\mathbf{b} - A_{j'} q_{j'l'}$ to $\mathbf{b} - A_{j'} Y_{j'l'}$ and the perturbation of fractions from \mathbf{q} to $\mathbf{q} + \Delta$, where Δ is a zero-mean random vector.

Let us re-formulate $P[\mathbf{b} - A_{j'} Y_{j'l'}, \mathbf{q} + \Delta, \mathbf{z}]$ based on how much we need to correct $\bar{\mathbf{x}}$:

$$\begin{aligned} (P[\mathbf{b} - A_{j'} Y_{j'l'}, \mathbf{q} + \Delta, \mathbf{z}]) \quad & \max_{\mathbf{y}} \quad \sum_{j,l} f_{jl}(q_{jl} + \Delta_{jl})(\bar{x}_{jl} - y_{jl}) \\ & \text{s.t.} \quad \sum_{j,l} a_{ij}(q_{jl} + \Delta_{jl})(\bar{x}_{jl} - y_{jl}) \leq b_i - a_{ij} Y_{j'l'} \quad \forall i \\ & \quad \quad \quad \sum_l (\bar{x}_{jl} - y_{jl}) \leq z_j \quad \forall j \\ & \quad \quad \quad \bar{\mathbf{x}} - \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

The new formulation uses decision variables \mathbf{y} , which may be negative, and correspond to how much we movement there is from the initial solution $\bar{\mathbf{x}}$ to the new one.

Let us denote the resource-slack variables of $P[\mathbf{b} - A_{j'}q_{j'l'}, \mathbf{q}, \mathbf{z}]$ by $(s_i \geq 0 : i \in [d])$, i.e., $\sum_{j,l} a_{ij}q_{jl}\bar{x}_{jl} + s_i = b_i - a_{ij'}q_{j'l'}$. Similarly, let us denote the demand-slack variables by $(u_j \geq 0 : j \in [n])$, i.e., $\sum_l \bar{x}_{jl} + u_j = z_j$. Using the slack variables, the problem simplifies to

$$\begin{aligned}
P[\mathbf{b} - A_{j'}Y_{j'l'}, \mathbf{q} + \Delta, \mathbf{z}] = \sum_{j,l} f_{jl}(q_{jl} + \Delta_{jl})\bar{x}_{jl} - \min_{\mathbf{y}} \quad & \sum_{j,l} f_{jl}(q_{jl} + \Delta_{jl})y_{jl} \\
\text{s.t.} \quad & \sum_{j,l} a_{ij}(q_{jl} + \Delta_{jl})y_{jl} \geq \beta_i \quad \forall i \\
& \sum_l y_{jl} \geq -u_j \quad \forall j \\
& \mathbf{y} \leq \bar{\mathbf{x}}, \quad (19)
\end{aligned}$$

where we defined $\beta_i := a_{ij'}(Y_{j'l'} - q_{j'l'}) - s_i + \sum_{j,l} a_{ij}\Delta_{jl}\bar{x}_{jl}$.

Observe that, since $\mathbb{E}[\Delta] = 0$, the first term outside the minimization, namely $\sum_{j,l} f_{jl}(q_{jl} + \Delta_{jl})\bar{x}_{jl}$, equals $\sum_{j,l} f_{jl}q_{jl}\bar{x}_{jl} = P[\mathbf{b} - A_{j'}q_{j'l'}, \mathbf{q}, \mathbf{z}]$ in expectation. The following result readily proves Lemma 4.

LEMMA 10 (Correction LP). *If we denote $\mathbf{q} = Q(t+1)$, then the Bellman Loss is bounded by $\mathbb{E}[L_B(t+1, \mathbf{b}, j', l')] \leq \mathbb{E}[P_C[Y_{j'l'}, \mathbf{q}, \Delta]]$, where $(P_C[Y_{j'l'}, \mathbf{q}, \Delta])$ is the minimization problem in Eq. (19). Furthermore, $\mathbb{E}[L_B(t+1, \mathbf{b}, j', l')] \leq 0$.*

PROOF. Recall that $\beta_i = a_{ij'}(Y_{j'l'} - q_{j'l'}) - s_i + \sum_{j,l} a_{ij}\Delta_{jl}\bar{x}_{jl}$ and observe that $\mathbb{E}[\beta_i] \leq 0$ for all i . We will find some deterministic values c_i such that the objective value of $P_C[Y_{j'l'}, \mathbf{q}, \Delta]$ is upper bounded by $\sum_i c_i\beta_i$, which proves the result.

We argue the upper bound on $(P_C[Y_{j'l'}, \mathbf{q}, \Delta])$ by bounding the optimal dual solution. The dual of $P_C[Y_{j'l'}, \mathbf{q}, \Delta]$ is

$$\max_{\mu, \lambda, \theta \geq 0} \left\{ \beta' \mu - u' \lambda - \sum_{j,l} \bar{x}_{jl} \theta_{jl} : (q_{jl} + \Delta_{jl})A'_{j'}\mu + \lambda_j - \theta_{jl} \leq f_{jl}(q_{jl} + \Delta_{jl}) \quad \forall j, l \right\}$$

This problem is the dual of a feasible and finite problem (see Eq. (19)), hence it has an optimal finite solution and we can bound $\mu_i \leq c_i$ for some deterministic values c_i . The objective value of this maximization problem is upper bounded by $\beta'c$, which proves the result. \square

C.2. Proof of Lemma 7

Recall that we wish to establish the following: if $\hat{\varphi}$ (used by ONLINE) has a solution with $x_{j'l'} = \max_l x_{jl} \gg 1$, then posting price f_{jl} is a satisfying action. To establish this, it remains to bound the difference between the LP $P_S[v_t, \mathbf{q}, \mathbf{z}]$ and its ‘‘perturbed’’ version $P_S[V_t, \mathbf{q} + \Delta\mathbf{q}, \mathbf{z} + \Delta\mathbf{z}]$. To that end, we first establish a bound on $v_t - V_t$; see item (i) in the discussion following Lemma 6.

LEMMA 11. *For fixed \mathbf{b} , denote $V_t = P[\mathbf{b}, Q(t), Z^t]$ and $v_t = P[\mathbf{b}, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]]$. If $t \geq c$, then, with probability at least $1 - c/t^2$, we have $v_t - V_t \geq -c\sqrt{t \log(t)}$. The constant c is independent of \mathbf{b} and depends on $(\mathbf{f}, F_1, \dots, F_n)$ only.*

PROOF. Set $\mathbf{q} = Q(t)$, $\mathbf{z} = Z^t$, $\Delta\mathbf{q} = \mathbb{E}[Q(t)] - Q(t)$, and $\Delta\mathbf{z} = \mathbb{E}[Z^t] - Z^t$. Take $\bar{\mathbf{x}}$ to be a solution of V_t and use a correction program analogous to Eq. (19) to conclude

$$\begin{aligned} v_t = V_t + \sum_{j,l} f_{jl} \Delta q_{jl} \bar{x}_{jl} - \min_{\mathbf{y}} \quad & \sum_{j,l} f_{jl} (q_{jl} + \Delta q_{jl}) y_{jl} \\ \text{s.t.} \quad & \sum_{j,l} a_{ij} (q_{jl} + \Delta q_{jl}) y_{jl} \geq \beta_i \quad \forall i \\ & \sum_l y_{jl} + t p_j \geq \sum_l \bar{x}_{jl} \quad \forall j \\ & \mathbf{y} \leq \bar{\mathbf{x}}, \end{aligned} \quad (20)$$

where $\beta_i = -s_i + \sum_{j,l} a_{ij} \Delta q_{jl} \bar{x}_{jl}$. We will argue an upper bound on the minimization problem by exhibiting a feasible solution.

Set $g(t) := \sqrt{\log(t)/t}$ and consider the solution $y_{jl} = \bar{x}_{jl} \frac{g(t)}{q_{jl} + \Delta q_{jl}}$. First recall that, by Lemma 5, $|\Delta q_{jl}| \leq g(t)$ with high probability. The objective value of this solution is $\sum_{j,l} f_{jl} \bar{x}_{jl} g(t)$, hence from Eq. (20) we get $v_t \geq V_t - 2 \sum_{j,l} f_{jl} g(t) \bar{x}_{jl}$. From here, using the fact that $\bar{\mathbf{x}}$ solves an LP with the constraint $\sum_l x_{jl} \leq Z_j^t$ for all j and that $Z_j^t \leq t$ a.s., we conclude the result by using that $\bar{x}_{jl} \leq t$.

We are left to check that our solution \mathbf{y} is feasible for the LP in Eq. (20). The first set of constraints is satisfied because $g(t) \geq \Delta q_{jl}$. The second set of constraints is satisfied since $\sum_l \bar{x}_{jl} \leq Z_j^t$ and $Z_j^t \leq t p_j + \sqrt{t \log(t)}$ w.h.p. Finally, the constraints $\mathbf{y} \leq \bar{\mathbf{x}}$ are satisfied since $g(t) \leq q_{jl} + \Delta q_{jl}$ for all t large enough. \square

PROOF OF LEMMA 7. Let us denote $\theta = (v, \mathbf{q}, \mathbf{z})$ and $\theta + \Delta\theta = (v + \Delta v, \mathbf{q} + \Delta\mathbf{q}, \mathbf{z} + \Delta\mathbf{z})$. Recall that \mathbf{b} and t are fixed throughout. The selection program for a fixed component (j', l') is given by

$$\begin{aligned} (P_S[\theta]) \quad & \max \quad x_{j'l'} \\ \text{s.t.} \quad & \sum_{j,l} f_{jl} q_{jl} x_{jl} \geq v \\ & \sum_{j,l} a_{ij} q_{jl} x_{jl} \leq b_i \quad \forall i \\ & \sum_l x_{jl} \leq z_j \quad \forall j \\ & \mathbf{x} \geq 0. \end{aligned}$$

If \bar{X} is the solution to $P[\mathbf{b}, Q(t), Z^t]$ (used by OFFLINE) and $\bar{\mathbf{x}}$ is the solution to $P[\mathbf{b}, \mathbb{E}[Q(t)], \mathbb{E}[Z^t]]$ (used by ONLINE), we want to prove $\bar{X}_{j'l'} \geq \bar{x}_{j'l'} - c\sqrt{t \log(t)}$. Equivalently, our aim is to prove the following:

$$P_S[\theta + \Delta\theta] \geq P_S[\theta] - c\sqrt{t \log(t)} \quad \text{w.p. } 1 - c/t^2.$$

We argue via Lagrangian relaxation. The Lagrangian of the selection problem with parameters $\theta + \Delta$ is given by

$$\begin{aligned} L(\mathbf{x}, \lambda; \theta + \Delta\theta) &= x_{j'l'} + \lambda_0 \left(\sum_{j,l} f_{jl} (q_{jl} + \Delta q_{jl}) x_{jl} - v - \Delta v \right) + \sum_i \lambda_i \left(b_i - \sum_{j,l} a_{ij} (q_{jl} + \Delta q_{jl}) x_{jl} \right) \\ &\quad + \sum_j \lambda_j \left(z_j + \Delta z_j - \sum_l x_{jl} \right) \\ &= L(\mathbf{x}, \lambda; \theta) + \lambda_0 \left(\sum_{j,l} f_{jl} \Delta q_{jl} x_{jl} - \Delta v \right) - \sum_i \lambda_i \sum_{j,l} a_{ij} \Delta q_{jl} x_{jl} + \sum_j \lambda_j \Delta z_j \end{aligned}$$

Define $D := \{\mathbf{x} : \mathbf{x} \geq 0, \|\mathbf{x}\|_\infty \leq t\}$. Observe that both $P_S[\theta]$ and $P_S[\theta + \Delta\theta]$ have solutions $\mathbf{x} \in D$. From Lemma 5 and Lemma 11 we have the following with probability $1 - c/t^2$:

$$|\Delta q_{ji} x_{jl}| \leq \sqrt{t \log(t)} \quad \forall \mathbf{x} \in D, \quad \Delta v \leq \sqrt{t \log(t)}, \quad \Delta z_j \geq \sqrt{t \log(t)}.$$

Let λ^* be the optimal dual variables of $P_S[\theta + \Delta\theta]$. We claim that there is a constant c such that $\|\lambda^*\|_\infty \leq c$. Assuming this claim, from the previous equation we get

$$L(\mathbf{x}, \lambda; \theta + \Delta\theta) \geq L(\mathbf{x}, \lambda; \theta) - c\sqrt{t \log(t)} \quad \forall \mathbf{x} \in D.$$

Using Strong Duality for the problem $P_S[\theta + \Delta\theta]$ we have

$$\begin{aligned} P_S[\theta + \Delta\theta] &= \max_{\mathbf{x} \geq 0} L(\mathbf{x}, \lambda^*; \theta + \Delta\theta) \\ &= \max_{\mathbf{x} \in D} L(\mathbf{x}, \lambda^*; \theta + \Delta\theta) \\ &\geq \max_{\mathbf{x} \in D} L(\mathbf{x}, \lambda^*; \theta) - c\sqrt{t \log(t)} \\ &\geq P_S[\theta] - c\sqrt{t \log(t)}. \end{aligned}$$

In the last step we used weak duality. Finally, to bound $\|\lambda^*\|_\infty \leq c$ we observe that the dual feasible region is defined by $\lambda \geq 0$ and the following set of inequalities, where δ is the Kronecker delta:

$$-f_{jl} q_{jl} \lambda_0 + q_{jl} \sum_i a_{ij} \lambda_i + \lambda_j \geq \delta_{j'l} \quad \forall j, l.$$

These inequalities are independent of (t, \mathbf{b}) , hence we can bound uniformly the extreme points. \square

Appendix D: Additional Details from Section 6 (Distribution-Agnostic Knapsack)

PROOF OF THEOREM 4. To apply Theorem 5, we first bound the measure of the exclusion sets \mathcal{B} and the ‘‘disagreement’’ sets \mathcal{Q} . Recall that $\mathcal{B}(t, b)$ is given in Lemma 8 and $\mathcal{Q}(t, b)$ is the event where \hat{U}^t is not a satisfying action.

Let $\sigma : [n] \rightarrow [n]$ be an ordering of $[n]$ w.r.t. the ratios $\bar{r}_j := \frac{r_j}{w_j}$ such that $\sigma_j = 1$ if j has the highest ratio. Similarly, let $\hat{\sigma}^t : [n] \rightarrow [n]$ be the ordering w.r.t. ratios $\bar{R}_j^t := R_j^t/w_j$.

Call E^t the event $\mathcal{B}(t, B^t) \cup \mathcal{Q}(t, B^t)$, then

$$\mathbb{P}[E^t] = \mathbb{P}[E^t, \sigma = \hat{\sigma}^t] + \mathbb{P}[E^t, \sigma \neq \hat{\sigma}^t] \leq \mathbb{P}[E^t, \sigma = \hat{\sigma}^t] + \mathbb{P}[\sigma \neq \hat{\sigma}^t].$$

Let N_j^t be the number of type- j samples observed by the beginning of period t . By definition, since we are given a sample of each type before the process starts, we have $N_j^t = Z_j^T - Z_j^t + 1$. Since the reward distribution is sub-Gaussian, it satisfies the Chernoff bound (Boucheron et al. 2013)

$$\mathbb{P}[R_j^t - r_j \geq x | N_j^t], \mathbb{P}[R_j^t - r_j \leq -x | N_j^t] \leq e^{-N_j^t x^2 / 2} \quad \forall x \in \mathbb{R}, \quad (21)$$

A union bound relying on Eq. (21) gives that

$$\mathbb{P}[\sigma \neq \hat{\sigma}^t | \mathcal{F}_t] \leq \mathbb{P}[\exists j \text{ s.t. } |\bar{r}_j - \bar{R}_j^t| \geq \delta/2 | \mathcal{F}_t] \leq 2 \sum_j e^{-N_j^t (w_j \delta)^2 / 8}.$$

The variable N_j^t , recall, is the number of type- j samples observed by the beginning of period t , hence $N_j^t - 1$ is a $\text{Bin}(T-t, p_j)$ random variable. It is a known fact that, given $\theta > 0$, $\mathbb{E}[e^{-\theta \text{Bin}(p, m)}] = (1 - p + pe^{-\theta})^m$, thus

$$\mathbb{P}[\sigma \neq \hat{\sigma}^t] = \mathbb{E}[\mathbb{P}[\sigma \neq \hat{\sigma}^t | \mathcal{F}_t]] \leq 2 \sum_j e^{-(w_j \delta)^2 / 8} (1 - p_j + p_j e^{-(w_j \delta)^2 / 8})^{T-t}.$$

Upper bounding by a geometric sum yields

$$\text{Regret}_2 := \sum_t \mathbb{P}[\sigma \neq \hat{\sigma}^t] \leq 2 \sum_j \frac{1}{p_j (e^{(w_j \delta)^2 / 8} - 1)} \leq 2 \sum_j \frac{8}{p_j (w_j \delta)^2}. \quad (22)$$

We are left to bound $\mathbb{P}[E^t, \sigma = \hat{\sigma}^t]$. Let us assume w.l.o.g. that the indexes are ordered so that $\bar{r}_1 \geq \bar{r}_2 \geq \dots \geq \bar{r}_n$. The optimal solution of $(P[B^t, \mathbf{r}, Z^t])$, i.e., OFFLINE's problem, is to sort the items and accept starting from $j = 1$, without exceeding the capacity B^t or the number of arrivals Z_j^t . Mathematically, the optimal solution X^{*t} to $(P[B^t, \mathbf{r}, Z^t])$ is

$$X_{1\mathbf{a}}^{*t} = \min \left\{ Z_1^t, \frac{B^t}{w_1} \right\}, \quad X_{j\mathbf{a}}^{*t} = \min \left\{ Z_j^t, \frac{B^t - \sum_{i < j} w_i X_{i\mathbf{a}}^{*t}}{w_j} \right\} \quad j = 2, \dots, n.$$

For the proxy $(P[B^t, R^t, \mu(t)])$, the optimal solution has the same structure with Z_j^t replaced everywhere by $\mu_j(t)$.

Let $\xi^t = j$ and U be any action in $\text{argmax}\{X_{j,u}^t : u = \mathbf{a}, \mathbf{r}\}$. We study first the case $U = \mathbf{a}$. If $X_{j,\mathbf{a}}^{*t} \geq 1$ then $U = \mathbf{a}$ would be, by Lemma 8, a satisfying action. If it is not a satisfying action it must then be that $X_{j,\mathbf{a}}^{*t} < 1$ and since the algorithm chooses to accept it must be also that $X_{j,\mathbf{a}}^t \geq \mu_j(t)/2$.

Thus we obtain the following two conditions

$$X_{j,\mathbf{a}}^{*t} < 1 \Rightarrow \sum_{i < j} w_i Z_i^t \geq b \quad \text{and} \quad X_{j,\mathbf{a}}^t \geq \mu_j(t)/2 \Rightarrow \sum_{i < j} w_i \mu_i(t) + w_j \mu_j(t)/2 \leq b.$$

In the case $U = \mathbf{r}$, $X_{j,\mathbf{r}}^{*t} < 1$ and $X_{j,\mathbf{r}}^t \geq \mu_j(t)/2$ imply

$$\sum_{i \leq j} w_i Z_i^t \leq b \quad \text{and} \quad \sum_{i < j} w_i \mu_i(t) + w_j \mu_j(t)/2 \geq b.$$

In conclusion,

$$\mathbb{P}[E^t, \sigma = \hat{\sigma}^t] \leq \max \left\{ \mathbb{P} \left[\sum_{i \leq j} w_i (Z_i^t - \mu_i(t)) \geq \frac{w_j \mu_j(t)}{2} \right], \mathbb{P} \left[\sum_{i \leq j} w_i (Z_i^t - \mu_i(t)) \leq -\frac{w_j \mu_j(t)}{2} \right] \right\}.$$

These probabilities are bounded symmetrically using the method of averaged bounded differences (Dubhashi and Panconesi 2009, Theorem 5.3). Indeed, using the natural linear function

$f(\xi^1, \dots, \xi^t) = \sum_i w_i \sum_{l=1}^t \mathbb{1}_{\{\xi^l=i\}}$, the differences are bounded by $|\mathbb{E}[f|\mathcal{F}_i] - \mathbb{E}[f|\mathcal{F}_{i-1}]| \leq w_{\max}$, hence

$$\text{Regret}_1 := \sum_t \mathbb{P}[E^t, \sigma = \hat{\sigma}^t] \leq \sum_t \sum_j p_j \exp\left(-\frac{2(w_j \mu_j(t)/2)^2}{t w_{\max}^2}\right) \leq 2 \sum_j \frac{(w_{\max}/w_j)^2}{p_j}.$$

Together with Eq. (22), we have the desired bound. \square

Appendix E: Connections to Information Relaxations

Our work is related to the information-relaxation framework developed in (Brown et al. 2010, Balseiro and Brown 2019). The information-relaxation framework is a fairly general way to endow OFFLINE with additional information, but at the same time forcing him to pay a penalty for using this information. The dual problem (with the penalties) is an upper bound on the performance of the best online policy.

The main distinctions with our approach are:

1. Information Relaxation requires to identify OFFLINE's filtration and penalties to build a proxy for OFFLINE's value function. This proxy can then be used to assess the performance of specific online policies.

The proxy that is developed—as the true OFFLINE value in our framework—may be difficult to compute. To overcome this difficulty, (Balseiro and Brown 2019) proposes an approximation through which penalties can be computed and hence an upper bound can be obtained.

2. Our framework requires, as well, identifying a suitable information structure (a filtration) and a relaxation φ . Because we allow for a Bellman Loss, we can develop φ $\hat{\varphi}$ that are computationally tractable. In most cases, a linear program. The framework explicitly then provides a mechanism, the RABBI algorithm, to derive a good online policy.

There is also an explicit mathematical connection. To state it, we first present a weaker version of our Bellman Inequalities, called thus because it is easier to find an object φ under this definition. Recall that, for a given non-anticipatory policy π , we denote v_π^{on} the expected value. Observe that the distinction with Definition 2 is in the initial ordering condition; we now require ϕ to upper bound the online value instead of the best offline.

DEFINITION 6 (WEAK BELLMAN INEQUALITIES). The sequence of r.v. $\{\varphi(t, s)\}_{t \in T, s \in \mathcal{S}}$ satisfies the Weak Bellman Inequalities w.r.t. filtration \mathcal{G} and events $\mathcal{B}(t, s) \subseteq \Omega$ if $\varphi(t, s)$ is \mathcal{G}_t -measurable for all t, s and the following holds:

1. Initial ordering: $\max_\pi v_\pi^{\text{on}} \leq \mathbb{E}[\varphi(T, S^T | \mathcal{G}_T)]$, where S^T is the initial state.
2. Monotonicity: $\forall s \in \mathcal{S}, t \in [T], \omega \notin \mathcal{B}(t, s)$,

$$\varphi(t, s | \mathcal{G}_t) \leq \max_{u \in \mathcal{U}} \{\mathcal{R}(s, \xi^t, u) + \mathbb{E}[\varphi(t-1, \mathcal{T}(s, \xi^t, u) | \mathcal{G}_{t-1}) | \mathcal{G}_t]\}. \quad (23)$$

In Proposition 2.1 in (Balseiro and Brown 2019) it is shown that if φ is some function that satisfies the Bellman equation for OFFLINE with the penalized immediate rewards function, then, in particular, it satisfies the initial ordering above. Since such φ satisfies, by construction, the Bellman inequality the following is an immediate corollary.

PROPOSITION 4 (**Proposition 2.1 in (Balseiro and Brown 2019)**). *Given feasible penalties z_t , the penalized value function satisfies Definition 2 with exclusion sets $\mathcal{B}(t, s) = \emptyset$.*

Our framework is a structured approach for building a computationally tractable φ , and deriving an online policy is bounded regret, without pre-computing penalties.

Appendix F: Parameter for the Pricing Instance

	Type j																			
Resource i	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1	1	0	0	1	0	1	0	0	0	0	1	0	1	1	0	1	1	0	0	1
2	0	0	0	0	1	1	1	1	0	0	0	0	0	0	0	1	0	1	0	0
3	0	1	0	1	0	1	1	1	0	0	1	1	1	1	1	0	1	1	1	1
4	1	0	1	1	1	0	0	1	0	0	0	0	0	0	1	1	0	0	1	1
5	1	1	1	1	0	0	0	1	1	0	0	0	0	0	0	0	1	1	0	0
6	0	0	0	0	1	1	1	1	1	1	0	0	0	0	1	0	1	1	0	1
7	0	0	0	0	1	1	1	1	1	0	0	1	1	0	0	0	1	1	0	1
8	0	1	0	0	0	0	1	1	1	1	1	0	0	0	1	1	1	1	0	1
9	1	0	0	0	0	0	1	1	0	1	1	1	1	1	0	0	0	0	0	0
10	1	0	0	0	0	0	0	0	0	0	1	0	1	1	0	1	0	0	1	0
11	1	0	1	1	0	0	0	1	0	0	1	1	1	1	0	0	1	0	0	1
12	1	0	1	0	1	1	0	0	1	0	0	1	1	0	1	0	1	0	0	0
13	1	1	1	0	0	0	1	1	1	0	1	1	0	0	1	1	0	0	0	0
14	1	0	1	0	0	0	1	0	0	1	1	1	0	0	1	0	1	1	0	1
15	1	1	0	1	1	0	1	0	0	1	0	1	1	1	0	0	1	1	0	0
16	0	1	0	1	0	0	1	0	0	1	0	1	1	1	0	0	0	1	1	0
17	1	0	1	1	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0
18	0	1	1	1	1	1	0	0	0	0	1	1	0	0	1	1	1	0	0	0
19	0	1	1	1	1	1	1	1	0	0	1	0	0	1	1	0	0	0	1	0
20	0	0	0	0	1	1	0	0	1	0	1	0	1	1	0	1	0	1	1	1
21	0	1	0	1	0	0	1	0	1	1	0	1	1	1	1	1	0	0	1	0
22	1	1	1	0	1	1	1	0	0	1	0	0	0	1	0	0	1	1	1	0
23	0	0	0	0	1	1	0	0	1	1	1	0	1	0	1	1	1	0	0	0
24	0	1	1	1	0	1	1	1	1	0	0	1	1	1	1	1	1	1	1	1
25	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	1	1	0
p_j	0.094	0.047	0.011	0.047	0.082	0.011	0.011	0.058	0.105	0.07	0.094	0.011	0.011	0.058	0.023	0.07	0.058	0.058	0.07	0.011
f_{j1}	18	4	5	13	4	3	20	15	17	16	20	16	16	13	8	5	10	13	20	16
f_{j2}	4	2	3	6	3	2	16	14	14	2	20	12	2	3	4	4	7	7	6	3
f_{j3}	1	1	2	2	2	1	7	6	6	1	1	1	1	1	3	3	1	1	2	2
q_{j1}	0.108	0.116	0.025	0.062	0.162	0.069	0.305	0.169	0.016	0.129	0.197	0.496	0.009	0.114	0.023	0.171	0.056	0.104	0.202	0.22
q_{j2}	0.329	0.21	0.495	0.233	0.229	0.223	0.458	0.33	0.191	0.215	0.579	0.966	0.046	0.154	0.105	0.648	0.291	0.137	0.618	0.27
q_{j3}	0.408	0.335	0.585	0.619	0.396	0.281	0.764	0.44	0.452	0.563	0.62	0.993	0.269	0.682	0.26	0.852	0.723	0.993	0.802	0.588