ADAPTIVE MAXIMIZATION OF SOCIAL WELFARE

NICOLÒ CESA-BIANCHI

Dept. of Computer Science, Università degli Studi di Milano and Politecnico di Milano

ROBERTO COLOMBONI

Dept. of Computer Science, Università degli Studi di Milano and Politecnico di Milano

MAXIMILIAN KASY

Department of Economics, University of Oxford

We consider the problem of repeatedly choosing policies to maximize social welfare. Welfare is a weighted sum of private utility and public revenue. Earlier outcomes inform later policies. Utility is not observed, but indirectly inferred. Response functions are learned through experimentation.

We derive a lower bound on regret, and a matching adversarial upper bound for a variant of the Exp3 algorithm. Cumulative regret grows at a rate of $T^{2/3}$. This implies that (i) welfare maximization is harder than the multiarmed bandit problem (with a rate of $T^{1/2}$ for finite policy sets), and (ii) our algorithm achieves the optimal rate. For the stochastic setting, if social welfare is concave, we can achieve a rate of $T^{1/2}$ (for continuous policy sets), using a dyadic search algorithm.

We analyze an extension to nonlinear income taxation, and sketch an extension to commodity taxation. We compare our setting to monopoly pricing (which is easier), and price setting for bilateral trade (which is harder).

KEYWORDS: Multiarmed bandits, optimal taxation, social welfare, adversarial learning.

1. INTRODUCTION

CONSIDER A POLICYMAKER who aims to maximize social welfare, defined as a weighted sum of utility across individuals. The policymaker can choose a policy parameter such as a sales tax rate, an unemployment benefit level, a health-insurance copay rate, etc. The policymaker does not directly observe the welfare resulting from their policy choices. They do, however, observe behavioral outcomes such as consumption of the taxed good, labor market participation, or health care expenditures. They can revise their policy choices over time in light of observed outcomes. How should such a policymaker act? To address this question, we bring together insights from welfare economics (in particular optimal taxation, Ramsey (1927), Mirrlees (1971), Baily (1978), Saez (2001), Chetty (2009)) with insights from machine learning (in particular online learning and multiarmed bandits, see Slivkins (2019), Lattimore and Szepesvári (2020) for recent reviews, and Thompson (1933), Lai and Robbins (1985) for classic contributions).

In our baseline model, individuals arrive sequentially and make a single binary decision. In each period, the policymaker chooses a tax rate that applies to this binary decision,

© 2025 The Authors. Econometrica published by John Wiley & Sons Ltd on behalf of The Econometric Society. Maximilian Kasy is the corresponding author on this paper. This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

Nicolò Cesa-Bianchi: nicolo.cesa-bianchi@unimi.it

Roberto Colomboni: roberto.colomboni@polimi.it

Maximilian Kasy: maximilian.kasy@economics.ox.ac.uk

NCB and RC were supported by the MUR PRIN grant 2022EKNE5K (Learning in Markets and Society), by the FAIR (Future Artificial Intelligence Research) project, funded by the NextGenerationEU program within the PNRR-PE-AI scheme, and by the EU Horizon research and innovation action under grant agreement 101120237, project ELIAS (European Lighthouse of AI for Sustainability). Maximilian Kasy was supported by the Alfred P. Sloan Foundation, under the grant "Social foundations for statistics and machine learning."

and then observes the individual's response. They do not observe the individual's private utility. Social welfare is given by a weighted sum of private utility and public revenue. Later, we extend our model to nonlinear income taxation, where welfare weights vary as a function of individual earnings capacity, and sketch an extension to commodity taxation, where individual decisions involve a continuous consumption vector.

Our goal is to give guidance to the policymaker. We propose algorithms to maximize cumulative social welfare, and we provide (adversarial and stochastic) guarantees for the performance of these algorithms. In doing so, we also show that welfare maximization is a harder learning problem than reward maximization in the multiarmed bandit setting. Private utility in our baseline model is equal to consumer surplus, which is given by the integral of demand. In order to learn this integral, we need to learn demand for counterfactual, suboptimal tax rates. This drives the difficulty of the learning problem.

Why Welfare, Why Adversarial Guarantees?. Our algorithms are designed to maximize social welfare, which is not directly observable, rather than maximizing outcomes that are directly observable. The definition of social welfare as an aggregation of individual utilities is at the heart of welfare economics in general, and of optimal tax theory in particular. The distinction between utility and observable outcomes is important in practice. To illustrate, consider the example of a policymaker who chooses the level of unemployment benefits, where the observable outcome is employment. The policymaker could use an algorithm that adaptively maximizes employment. The problem with this approach is that employment might be maximized by making the unemployed as miserable as possible. This is not normatively appealing. Such an algorithm would minimize the utility of the unemployed, rather than maximizing social welfare. Similar examples can be given for many domains of public policy, including health, education, and criminal justice. In contrast to observable outcomes such as employment, welfare is improved by increasing the choice sets of those affected, not by reducing these choice sets.

Our theoretical analysis provides not only stochastic but also adversarial guarantees, which hold for arbitrary sequences of preference parameters. Adversarial guarantees for algorithms promise robustness against deviations from the assumption that heterogeneity is independently and identically distributed over time. Possible deviations from this assumption include autocorrelation, trends, heteroskedasticity, more general nonstationarity, and other concerns of time-series econometrics. In the employment example, such deviations might for instance be due to the business cycle. One might fear that adversarial robustness is achieved at the price of worsened performance for the i.i.d. setting, relative to less robust algorithms. That this is not the case follows from our theoretical characterizations.

Lower and Upper Bounds on Regret. Our main theorems provide lower and upper bounds on cumulative regret. Cumulative regret is defined as the difference in welfare between the chosen sequence of policies and the best possible constant policy. We consider both stochastic and adversarial regret. A lower bound on stochastic regret satisfies that, for any algorithm, there exists some stationary distribution of preference parameters for which the algorithm has to suffer at least a certain amount of regret. An upper bound on adversarial regret has to hold for a given algorithm and any sequence of preference parameters.

For a given algorithm, stochastic regret, averaged over i.i.d. sequences of preference parameters, is always less or equal than adversarial regret, for the worst-case sequence. A lower bound on stochastic regret (for any algorithm) therefore implies a corresponding lower bound on adversarial regret, and an upper bound on adversarial regret (for a given algorithm) immediately implies an upper bound on stochastic regret. When an adversarial upper bound coincides with a stochastic lower bound, in terms of rates of regret, it follows that the proposed algorithm is rate efficient, for both stochastic and adversarial regret. It follows, furthermore, that the bounds are sharp.

A Lower Bound on Stochastic Regret. We first prove a stochastic (and thus also adversarial) lower bound on regret, for any possible algorithm in the welfare maximization problem. Our proof of this bound constructs a family of possible distributions for preferences. This family is such that there are two candidate policies, which are potentially optimal. The difference in welfare between these two policies depends on the integral of demand over intermediate policy values. In order to learn which of the two candidate policies is optimal, we need to learn behavioral responses for intermediate policies, which are strictly suboptimal. Because of the need to probe these suboptimal policies sufficiently often, we obtain a lower bound on regret, which grows at a rate of $T^{2/3}$, even if we restrict our attention to settings with finite, known support for preference parameters and policies. This rate is worse than the worst-case rate for bandits of $T^{1/2}$.

A Matching Upper Bound on Adversarial Regret for Modified Exp3. We next propose an algorithm for the adaptive maximization of social welfare. Our algorithm is a modification of the Exp3 algorithm (Auer, Cesa-Bianchi, Freund, and Schapire (2002)). Exp3 is based on an unbiased estimate of cumulative welfare for each policy. The probability of choosing a given policy is proportional to the exponential of this estimate of cumulative welfare, times some rate parameter. Relative to Exp3, we require two modifications for our setting. First, we need to discretize the continuous policy space. Second, and more interestingly, we need additional exploration of counterfactual policies, including some policies that are clearly suboptimal, in order to learn welfare for the policies, which are contenders for the optimum. This need for additional exploration again arises because of the dependence of welfare on the integral of demand over counterfactual policy choices. For our modified Exp3 algorithm, we prove an adversarial (and thus also stochastic) upper bound on regret. We show that, for an appropriate choice of tuning parameters, worst-case cumulative regret over all possible sequences of preference parameters grows at a rate of $T^{2/3}$, up to a logarithmic term. The algorithm thus achieves the best possible rate. Since the rates for our stochastic lower and adversarial upper bound coincide, up to a logarithmic term, we have a complete characterization of learning rates for the welfare maximization problem.

Improved Stochastic Bounds for Concave Social Welfare. The proof of our lower bound on regret is based on the construction of a distribution of preferences which delivers a nonconcave social welfare function. If we restrict attention to the stochastic setting, where preferences are i.i.d. over time, and if we assume that social welfare is concave, then we can improve upon this bound on regret. We prove a lower bound on stochastic regret, under the assumption of concavity, which grows at the rate of $T^{1/2}$. We then propose a dyadic search algorithm, which achieves this rate, up to logarithmic terms. This dyadic search algorithm maintains an "active interval," containing the optimal policy with high probability, which is narrowed down over time. Only policies within the active interval are sampled.

Extensions to Nonlinear Income Taxation and to Commodity Taxation. Our discussion up to this point focuses on a minimal, stylized case of an optimal tax problem, where individual actions are binary, and the policy imposes a tax on this binary action. Our arguments generalize, however, to more complicated and practically relevant settings. This includes optimal nonlinear income taxation, as in Mirrlees (1971) and Saez (2001), and commodity taxation for a bundle of goods, as in Ramsey (1927). For nonlinear income taxation, different tax rates apply at different income levels, and welfare weights depend on individual earnings capacity. In Section 5, we discuss an extension of our tempered Exp3 algorithm to nonlinear income taxation, and characterize its regret. For commodity taxation, different tax rates apply to different goods, and consumption decisions are continuous vectors. In Section 6, we sketch an extension of our algorithm to commodity taxation, but leave its characterization for future research.

Roadmap. The rest of this paper proceeds as follows. We conclude this introduction with a discussion of some related work and relevant references. Section 2 introduces our setup, formally defines the adversarial and stochastic settings, and compares our setup to related learning problems. Section 3 provides lower and upper bounds on regret in the adversarial and stochastic settings. Section 4 restricts attention to the stochastic setting with concave social welfare, and provides improved regret bounds for this setting. Section 5 discusses an extension of our baseline model to nonlinear income taxation. Section 6 sketches another extension of our baseline model to commodity taxation. Section 7 concludes, and discusses some possible applications of our algorithm, as well as an alternative Bayesian approach to adaptive welfare maximization. The proofs of Theorem 1 and Theorem 2 can be found in Appendix A. The proofs of our remaining theorems and proofs of technical lemmas are discussed in the Online Supplement (Cesa-Bianchi, Colomboni, and Kasy (2025)).

1.1. Background and Literature

To put our work in context, it is useful to contrast our framework with the standard approach in public finance and optimal tax theory, and with the frameworks considered in machine learning and the multiarmed bandit literature.

Optimal Taxation. Optimal tax theory, and optimal policy theory more generally, is concerned with the maximization of social welfare, where social welfare is understood as a (weighted) sum of subjective utility across individuals (Ramsey (1927), Mirrlees (1971), Baily (1978), Saez (2001), Chetty (2009)). A key tradeoff in such models is between, first, redistribution to those with higher welfare weights, and second, the efficiency cost of behavioral responses to tax increases. Such behavioral responses might reduce the tax base.

Optimal tax problems are defined by normative parameters (such as welfare weights for different individuals), as well as empirical parameters (such as the elasticity of the tax base with respect to tax rates). The typical approach in public finance uses historical or experimental variation to estimate the relevant empirical parameters (causal effects, elasticities). These estimated parameters are then plugged into formulas for optimal policy choice, which are derived from theoretical models. The implied optimal policies are finally implemented, without further experimental variation.

Multiarmed Bandits. The standard approach of public finance, which separates elasticity estimation from policy choice, contrasts with the adaptive approach that characterizes decision-making in many branches of AI, including online learning, multi-armed bandits, and reinforcement learning. Multiarmed bandit algorithms, in particular, trade off exploration and exploitation over time (Slivkins (2019), Lattimore and Szepesvári (2020)).

1076

Exploration here refers to the acquisition of information for better future policy decisions, while exploitation refers to the use of currently available information for optimal policy decisions at the present moment. The goal of bandit algorithms is to maximize a stream of rewards, which requires an optimal balance between exploration and exploitation. Bandit algorithms for the stochastic setting are characterized by optimism in the face of uncertainty: Policies with uncertain payoff should be tried until their expected payoff is clearly suboptimal.

Bandit algorithms (and similarly, adaptive experimental designs for informing policy choice, as in Russo (2020), Kasy and Sautmann (2021)) are not directly applicable to social welfare maximization problems, such as those of optimal tax theory. The reason is that bandit algorithms maximize a stream of observed rewards. By contrast, social welfare as conceived in welfare economics is based on unobserved subjective utility.

Adversarial Decision-Making. Adversarial models for sequential decision-making find their roots in repeated game theory (Hannan (1957)), while related settings were independently studied in information theory (Cover (1965)) and computer science (Vovk (1990), Littlestone and Warmuth (1994)). Regret minimization, also in a bandit setting, was investigated as a tool to prove convergence of uncoupled dynamics to equilibria in *N*-person games (Hart and Mas-Colell (2000, 2001))—the exponential weighting scheme used by Exp3 is also known as *smooth fictitious play* in the game-theoretic literature (Fudenberg and Levine (1995)). Recent works (Seldin and Slivkins (2014), Zimmert and Seldin (2021)) show that simple variants of Exp3 simultaneously achieve essentially optimal regret bounds in adversarial, stochastic, and contaminated settings, without prior knowledge of the actual regime. This suggests that algorithms designed for adversarial environments can behave well in more benign settings, whereas the opposite is provably not true.

Bandit Approaches for Economic Problems. Bandit-type approaches have been applied to a number of other economic and financial scenarios in the literature where rewards are observable. These include monopoly pricing (Kleinberg and Leighton (2003)) (see also the survey by den Boer (2015)), second-price auctions (Weed, Perchet, and Rigollet (2016)), first-price auctions (Han, Zhou, and Weissman (2020), Han, Zhou, Flores, Ordentlich, and Weissman (2020), Achddou, Cappé, and Garivier (2021)); see also Kolumbus and Nisan (2022), Feng, Podimata, and Syrgkanis (2018), Feng, Guruganesh, Liaw, Mehta, and Sethi (2021), and combinatorial auctions (Daskalakis and Syrgkanis (2022)). Bandit-type approaches have also been applied to some settings where rewards are not directly observable, including bilateral trading (Cesa-Bianchi, Cesari, Colomboni, Fusco, and Leonardi (2024a, 2024b), and the newsvendor problem (Lugosi, Markakis, and Neu (2023)).

Bandit algorithms are widely used in online advertising and recommendation. Online learning methods are successfully used for tuning the bids made by autobidders (a service provided by advertising platforms) (Lucier, Pattathil, Slivkins, and Zhang (2024)). While these algorithms are analyzed in adversarial environments, the extent to which they are deployed in commercial products remains unclear.

2. Setup

At each time i = 1, 2, ..., T, one individual arrives who is characterized by an unknown willingness to pay $v_i \in [0, 1]$. This individual is exposed to a tax rate x_i , and makes a binary decision $y_i = \mathbf{1}(x_i \le v_i)$. The implied public revenue is $x_i \cdot y_i$. The implied private welfare is

 $\max(v_i - x_i, 0)$. We define social welfare as a weighted sum of public revenue and private welfare, with a weight $\lambda \in (0, 1)$ for the latter. Social welfare for time period *i* is therefore given by

$$U_i(x_i) = \underbrace{x_i \cdot \mathbf{1}(x_i \le v_i)}_{\text{Public revenue}} + \lambda \cdot \underbrace{\max(v_i - x_i, 0)}_{\text{Private welfare}}.$$
 (1)

After period *i*, we observe y_i and the tax rate x_i , but nothing else. In particular, we do not observe welfare $U_i(x_i)$.

We can rewrite social welfare $U_i(x)$ as follows. Denote $G_i(x) = \mathbf{1}(v_i \ge x)$, so that $y_i = G_i(x_i)$. This is the individual demand function. Then private welfare can be written as $\max(v_i - x, 0) = \int_x^1 G_i(x') dx'$. That is, private welfare is given by integrated demand.¹ This representation of private welfare implies

$$U_{i}(x) = \underbrace{x \cdot G_{i}(x)}_{\text{Public revenue}} + \lambda \cdot \underbrace{\int_{x}^{1} G_{i}(x') \, \mathrm{d}x'}_{\text{Private welfare}}.$$
(2)

We consider algorithms for the choice of x_i , which might depend on the observable history $(x_j, y_j)_{i=1}^{i-1}$, as well as possibly a randomization device.

Notation. For the *adversarial* setting, we will consider cumulative demand and welfare, denoted by blackboard bold letters, summing across j = 1, ..., i. In particular,

$$\mathbb{G}_i(x) = \sum_{j \leq i} G_i(x), \qquad \mathbb{U}_i(x) = \sum_{j \leq i} U_i(x), \qquad \mathbb{U}_i = \sum_{j \leq i} U_j(x_j).$$

 $\mathbb{G}_i(x)$ and $\mathbb{U}_i(x)$ are cumulative demand and welfare for a counterfactual, fixed policy x. \mathbb{U}_i , without an argument, is the cumulative welfare for the policies x_i actually chosen.

For the *stochastic* setting, we will analogously consider expected demand and expected welfare, denoted by boldface letters. The expectation is taken across some stationary distribution μ of v_i , where v_i is statistically independent of x_i , and of v_i for $j \neq i$. In particular,

$$\boldsymbol{G}(\boldsymbol{x}) = E[\boldsymbol{G}_i(\boldsymbol{x})], \qquad \boldsymbol{U}(\boldsymbol{x}) = E[\boldsymbol{U}_i(\boldsymbol{x})].$$

2.1. Regret

The Adversarial Case. Following the literature, we consider regret for both the adversarial and the stochastic setting. In the adversarial setting, we allow for arbitrary sequences of willingness to pay, $\{v_i\}_{i=1}^T$. We compare the expected performance of any given algorithm for choosing $\{x_i\}_{i=1}^T$ to the performance of the best possible constant policy x. This comparison yields cumulative expected regret, which is given by

$$\mathcal{R}_T(\lbrace v_i \rbrace_{i=1}^T) = \sup_{x} E[\mathbb{U}_T(x) - \mathbb{U}_T | \lbrace v_i \rbrace_{i=1}^T].$$
(3)

The expectation in this expression is taken over any possible randomness in the tax rates x_i chosen by the algorithm; there is no other source of randomness.

¹This reflects the absence of income effects in our model, which implies that private utility, consumer surplus, compensating variation, and equivalent variation all coincide.

The Stochastic Case. We also consider the stochastic setting. In this setting, we add structure by assuming that the v_i are i.i.d. draws from some distribution μ on [0, 1], with implied demand function $G(x) = P(v_i \ge x)$. This demand function is identified by the regression

$$G(x) = E[y_i | x_i = x].$$

The expectation in this expression is taken over the distribution of v_i , which is presumed to be statistically independent of the tax rate x_i . Expected welfare for this distribution of v_i is given by

$$U(x) = x \cdot G(x) + \lambda \int_{x}^{1} G(x') \, \mathrm{d}x'.$$

Cumulative expected regret in the stochastic case equals

$$\mathcal{R}_T(G) = \sup_x E\left[\mathbb{U}_T(x) - \mathbb{U}_T\right] = T \cdot \sup_x U(x) - E\left[\sum_{i \le T} U(x_i)\right].$$
(4)

The expectation in this expression is taken over both any possible randomness in the tax rates x_i , and the i.i.d. draws v_i .

2.2. Comparison to Related Learning Problems

Before proceeding with our analysis of regret, we take a step back, and compare our learning problem to two related problems that have received some attention in the literature. The first of these is the adaptive *monopoly pricing* problem; see, for instance, Kleinberg and Leighton (2003). This problem is equivalent to our setting when we set $\lambda = 0$, interpret x as a price, and U_i^{MP} as monopolist profits (neglecting production costs):

$$U_i^{\text{MP}}(x) = x_i \cdot \mathbf{1}(x_i \le v_i) = \underbrace{x \cdot G_i(x)}_{\text{Monopolist revenue}}.$$
(5)

As in our adaptive taxation setting, the feedback received at the end of period *i* is

$$y_i = G_i(x_i) = \mathbf{1}(x_i \le v_i).$$

Another related problem is price setting for *bilateral trade*; see, for instance, Cesa-Bianchi et al. (2024a). In this problem, welfare $U_i^{\text{BT}}(x)$ is given by the sum of seller and buyer welfare. Trade happens if and only if both sides agree to transact at the proposed price. Buyer willingness to pay is given by v_i^b , while the seller is willing to trade at prices above v_i^s :

$$U_i^{\text{BT}}(x) = \mathbf{1} \left(v_i^b \ge x \right) \cdot \max \left(x - v_i^s, 0 \right) + \mathbf{1} \left(v_i^s \le x \right) \cdot \max \left(v_i^b - x, 0 \right)$$
$$= G_i^b(x) \cdot \underbrace{\int_0^x G_i^s(x') \, \mathrm{d}x'}_{\text{Seller welfare}} + G_i^s(x) \cdot \underbrace{\int_x^1 G_i^b(x') \, \mathrm{d}x'}_{\text{Buyer welfare}}.$$
(6)

Feedback in this case is a little richer: We observe both whether the buyer *b* would have accepted the posted price, and whether the seller would have accepted this price,

$$y_i^b = G_i^b(x_i) = \mathbf{1}(x_i \le v_i^b)$$
 and $y_i^s = G_i^s(x_i) = \mathbf{1}(x_i \ge v_i^s).$

Model	Policy Space		Objective Function	
	Discrete	Continuous	Pointwise	One-Sided Lipschitz
Monopoly price setting	$T^{1/2}_{7}$	$T^{2/3}_{2/3}$	Yes	Yes
Bilateral trade	$T^{2/3}$ $T^{2/3}$	$T^{2/3}$	No No	Yes No

 TABLE I

 Regret rates for different learning problems.

Note: This table shows the efficient rates of regret for different learning problems. Rates are up to logarithmic terms, and apply to both the stochastic and the adversarial setting. Regret rates are shown for the discrete case, where the space of policies x is restricted to a finite set, and the continuous case, where x can take any value in [0, 1]. The columns on the right describe the properties of the objective function in each problem, which drive the differences in regret rates. Rates for the optimal taxation case are proven in this paper. Rates for the continuous monopoly price setting case are from Kleinberg and Leighton (2003); the discrete case reduces to a standard bandit problem. Rates for the continuous bilateral trade case are from Cesa-Bianchi et al. (2024a); the discrete case can be deduced by adapting the arguments in the same paper (for the stochastic i.i.d. case with independent sellers' and buyers' valuations), or by adapting the techniques in Cesa-Bianchi et al. (2024b) (for the adversarial case, allowing the learner to use weakly budget balanced mechanisms).

Lipschitzness and Information Requirements. The difficulty of the learning problem in each of these models critically depends on (i) the Lipschitz properties of the welfare function, and (ii) the information required to evaluate welfare at a point.

We say that a generic welfare function $W : [0, 1] \to \mathbb{R}$ is one-sided Lipschitz if $W(x + \varepsilon) \le W(x) + \varepsilon$ for all $0 \le x \le 1$ and all $0 \le \varepsilon \le 1 - x$. One-sided Lipschitzness allows us to bound the approximation error of a learning algorithm operating on a finite subset of the set of policies. One-sided Lipschitzness is an intrinsic property of both the monopoly pricing and the optimal taxation problem; it is not an assumption that is additionally imposed. To see this for monopoly pricing, note that, for $\epsilon \ge 0$, $U_i^{\text{MP}}(x + \varepsilon) = (x + \varepsilon) \cdot \mathbf{1}(x + \varepsilon \le v_i) \le x \cdot \mathbf{1}(x \le v_i) + \varepsilon = U_i^{\text{MP}}(x) + \varepsilon$. For social welfare, $U_i(x) = (x_i + \varepsilon) \cdot \mathbf{1}(x_i + \varepsilon \le v_i) + \lambda \cdot \max(v_i - x_i - \varepsilon, 0) \le x \cdot \mathbf{1}(x \le v_i) + \varepsilon + \lambda \cdot \max(v_i - x_i, 0) = U_i(x) + \varepsilon$.

We say that learning $W(\cdot)$ requires only pointwise information if W(x) is a function of G(x), and does not depend on $G(\cdot)$ otherwise. Pointwise information allows us to avoid exploring policies that are clearly suboptimal, when we aim to learn the optimal policy.

Table I summarizes the Lipschitz properties and information requirements in each of the three models; the following justifies the claims made in Table I:

- 1. For monopoly pricing, welfare $U_i^{MP}(x)$ is one-sided Lipschitz and only depends on $G_i(x)$ pointwise.
- 2. For *optimal taxation*, welfare $U_i(x)$ is one-sided Lipschitz and depends on both $G_i(x)$ at the given x (pointwise), and on an integral of $G_i(x')$ for a range of values of x' (nonpointwise).
- 3. For *bilateral trade*, welfare $U_i^{\text{BT}}(x)$ is not one-sided Lipschitz and depends on both $G_i^b(x)$ and $G_i^s(x)$ (pointwise), as well as the integrals of $G_i^b(x')$ and $G_i^s(x')$ (non-pointwise).

These properties suggest a ranking in terms of the difficulty of the corresponding learning problems, and in particular in terms of the rates of divergence of cumulative regret: The information requirements of optimal taxation are stronger than those of monopoly pricing, but its continuity properties are more favorable than those of bilateral trade. This intuition is correct, as shown by Table I. The rates for monopoly pricing and for bilateral trade are known (or can be easily adapted) from the literature. In this paper, we prove corresponding rates for optimal taxation. In comparing optimal taxation and monopoly pricing to conventional multiarmed bandits, it is worth emphasizing that there are two distinct reasons for the slower rate of convergence. First, the continuous support of x, as opposed to a finite number of arms, which is shared by optimal taxation and monopoly pricing. Second, the requirement of additional exploration of suboptimal policies for the optimal tax problem. As shown in Table I, the continuous support alone is enough to slow down convergence, with no extra penalty for the additional exploration requirement, in terms of rates. If, however, we restrict our attention to a discrete set of feasible policies x, then monopoly pricing reduces to a multiarmed bandit problem, with a minimax regret rate of $T^{1/2}$. The optimal tax problem, by contrast, still has a rate of $T^{2/3}$, even if we restrict our attention to the case of finite known support for v and x, as shown by the proof of Theorem 1 below.

Hannan Consistency. The cumulative regret of any nonadaptive algorithm necessarily grows at a rate of T. This includes, in particular, randomized experiments where the policy is chosen uniformly at random, from a fixed policy set, in every period. Algorithms for which adversarial regret (and thus also stochastic regret) grows at a rate less than T, so that per period regret goes to 0 as T increases, are known as *Hannan consistent*. Nonadaptive algorithms are not Hannan consistent. Table I implies that Hannan consistent algorithms exist in all settings considered, with the exception of Bilateral trade and continuous policy spaces.

3. STOCHASTIC AND ADVERSARIAL REGRET BOUNDS

We now turn to our main theoretical results, lower and upper bounds on stochastic and adversarial regret for the problem of social welfare maximization. We first prove a lower bound on stochastic regret, which applies to any algorithm, and which immediately implies a lower bound on adversarial regret. We then introduce the algorithm Tempered Exp3 for Social Welfare. We show that, for an appropriate choice of tuning parameters, this algorithm achieves the rates of the lower bound on regret, up to a logarithmic term. Formal proofs of these bounds can be found in Appendix A.

3.1. Lower Bound

THEOREM 1—Lower Bound on Regret: Consider the setup of Section 2. There exists a constant C > 0 such that, for any randomized algorithm for the choice of $x_1, x_2, ...$ and any time horizon $T \in \mathbb{N}$, the following holds:

- 1. There exists a distribution μ on [0, 1] with associated demand function G for which the stochastic cumulative expected regret $\mathcal{R}_T(G)$ is at least $C \cdot T^{2/3}$.
- 2. There exists a sequence (v_1, \ldots, v_T) for which the adversarial cumulative expected regret $\mathcal{R}_T(\{v_i\}_{i=1}^T)$ is at least $C \cdot T^{2/3}$.

The proof of Theorem 1 can be found in Appendix A. The adversarial lower bound follows immediately from the stochastic lower bound, since worst-case regret (over possible sequences of v_i) is bounded below by average regret (over i.i.d. draws of v_i), for any distribution of v_i .

Sketch of Proof. To prove the stochastic lower bound, we construct a family of distributions $\{\mu^{\epsilon}\}_{\epsilon \in [-1,1]}$ for v_i , indexed by a parameter $\epsilon \in [-1, 1]$. The distributions in this family have four points of support, (1/4, 1/2, 3/4, 1). The probability of these points is given by

$$(a, (1+\epsilon)b, (1-\epsilon)b, 1-a-2b).$$



FIGURE 1.—Construction for proving the lower bound on regret. *Notes*: This figure illustrates our construction for proving the lower bound on regret. The relative social welfare of policies 1 and .25 depends on the sign of ϵ . The solid line corresponds to $\epsilon = -1$, the dashed line to $\epsilon = 1$. In order to distinguish between these two, we must learn demand in the intermediate interval [0.5, 0.75].

The values of a and b are chosen such that (i) the two middle points $\frac{1}{2}, \frac{3}{4}$ are far from optimal, for any value of ϵ , and (ii) learning which of the two end points $\binom{1}{4}, 1$ is optimal requires sampling from the middle.² For each $\epsilon \in [-1, 1]$, denote the demand function associated to μ^{ϵ} by G^{ϵ} , and the expected social welfare associated to G^{ϵ} by U^{ϵ} . Property (ii) holds because of the integral term $\int_{\frac{1}{4}}^{1} G^{\epsilon}(x') dx'$, which shows up in $U^{\epsilon}(1) - U^{\epsilon}(\frac{1}{4})$. This construction is illustrated in Figure 1. This figure shows plots of G^{ϵ} and of U^{ϵ} for $\lambda = 0.95$ and $\epsilon \in \{\pm 1\}$.

The difference in welfare $U^{\epsilon}(1) - U^{\epsilon}(1/4)$ of the two candidate optimal policies 1/4 and 1 depends on the sign of ϵ . In order not to suffer expected regret that grows as $|\epsilon| \cdot T$, any learning algorithm needs to sample policies from points that are informative about the sign of ϵ . The only points that are informative are those in the region (1/2, 3/4], where welfare is bounded away from optimal welfare.

More specifically, the learning algorithm has to sample on the order of $|\epsilon|^{-2}$ times from the region (1/2, 3/4], to be able to detect the sign of ϵ , incurring regret on the order of $|\epsilon|^{-2}$ in the process. Any learning algorithm therefore incurs regret on the order of $\min(|\epsilon|^{-2}, |\epsilon| \cdot T)$, which for $\epsilon \propto T^{-1/3}$), leads to the conclusion.

3.2. An Algorithm That Achieves the Lower Bound

We next introduce Algorithm 1, which allows us to essentially achieve the lower bound on regret, in terms of rates.

Conventional Exp3. Algorithm 1 is a modification of the Exp3 algorithm. Conventional Exp3 (Auer et al. (2002)) is designed to maximize the standard bandit objective, $\sum_{i \leq T} y_i$. Exp3 maintains an unbiased running estimate of the cumulative payoff of each arm k, calculated using inverse probability weighting, $\widehat{\mathbb{G}}_{i,k} = \sum_{j < i} y_i \cdot \frac{\mathbf{1}(k_i = k)}{p_{ik}}$. In period i, arm k is chosen with probability $p_{ik} = (1 - \gamma) \cdot \frac{\exp(\eta \cdot \widehat{\mathbb{G}}_{ik})}{\sum_{k'} \exp(\eta \cdot \widehat{\mathbb{G}}_{ik'})} + \frac{\gamma}{K+1}$, where η and γ are tuning parameters. p_{ik} is thus increasing in the estimated average performance $\frac{\widehat{\mathbb{G}}_{i,k}}{i}$ of arm k in prior periods. Because $\widehat{\mathbb{G}}_{i,k}$ is *not* normalized by the number of time periods i, more

²Specifically, $a := \frac{(1-\lambda)\cdot(136-99\cdot\lambda)}{2\cdot(4-3\cdot\lambda)\cdot(24-17\cdot\lambda)}$, and $b := \frac{1-\lambda}{2\cdot(24-17\cdot\lambda)}$. These two constants are strictly greater than zero, and satisfy $1 - a - 2 \cdot b > 0$.

Algorithm 1 Tempered Exp3 for Social Welfare.

Require: Tuning parameters K, γ and η .

- 1: Calculate evenly spaced grid-points $\tilde{x}_k = (k-1)/K$, and initialize $\widehat{\mathbb{G}}_{1k} = 0$ and $\widehat{\mathbb{U}}_{1k} = 0$ for $k = 1, \dots, K+1$.
- 2: **for** individual i = 1, 2, ..., T **do**
- 3: For all k = 1, 2, ..., K + 1, set

{Assignment probabilities}

$$p_{ik} = (1 - \gamma) \cdot \frac{\exp(\eta \cdot \mathbb{U}_{ik})}{\sum_{k'} \exp(\eta \cdot \widehat{\mathbb{U}}_{ik'})} + \frac{\gamma}{K+1}.$$
 (7)

- 4: Choose k_i at random according to the probability distribution $(p_{i,1}, \ldots, p_{i,K+1})$. Set $x_i = \tilde{x}_{k_i}$, and query y_i accordingly. For all $k = 1, 2, \dots, K + 1$ act
- 5: For all k = 1, 2, ..., K + 1, set

$$\widehat{\mathbb{G}}_{i+1,k} = \widehat{\mathbb{G}}_{i,k} + y_i \cdot \frac{\mathbf{1}(k_i = k)}{p_{ik}}.$$
(8)

6: For all k = 1, 2, ..., K + 1, set

{Estimated welfare}

$$\widehat{\mathbb{U}}_{i+1,k} = \widetilde{x}_k \cdot \widehat{\mathbb{G}}_{i+1,k} + \frac{\lambda}{K} \cdot \sum_{k' > k} \widehat{\mathbb{G}}_{i+1,k'}.$$
(9)

7: end for

weight is given to the best performing arms over time, as estimation uncertainty for average performance decreases. In both these aspects, Exp3 is similar to the popular Upper Confidence Bound algorithm (UCB) for stochastic bandit problems (Lai (1987), Agrawal (1995), Auer, Cesa-Bianchi, and Fischer (2002)). In contrast to UCB, Exp3 is a randomized algorithm. Randomization is required for adversarial performance guarantees. This is analogous to the necessity of mixed strategies for zero-sum games.

Modifications Relative to Conventional Exp3. Relative to this algorithm, we require three modifications. First, we discretize the continuous support [0, 1] of x, restricting attention to the grid of policy values $\tilde{x}_k = (k-1)/K$. Second, since welfare $U_i(x)$ is not directly observed for the chosen policy x, we need to estimate it indirectly. In particular, we first form an estimate $\widehat{\mathbb{G}}_{ik}$ of cumulative demand for each of the policy values \tilde{x}_k , using inverse probability weighting. We then use this estimated demand, interpolated using a stepfunction, to form estimates of cumulative social welfare, $\widehat{\mathbb{U}}_{ik} = \tilde{x}_k \cdot \widehat{\mathbb{G}}_{ik} + \frac{\lambda}{K} \cdot \sum_{k'>k} \widehat{\mathbb{G}}_{ik'}$. Third, we require additional exploration, relative to Exp3. Since social welfare depends on demand for counterfactual policy choices, we need to explore policies that are away from the optimum, in order to learn the relative welfare of approximately optimal policy choices. The mixing weight γ , which determines the share of policies sampled from the uniform distribution, needs to be larger relative to conventional Exp3, to ensure sufficient exploration away from the optimum.

THEOREM 2—Adversarial Upper Bound on Regret of Tempered Exp3 for Social Welfare: Consider the setup of Section 2, and Algorithm 1. Assume that $(K + 1)\eta < \gamma$. Then for any sequence (v_1, \ldots, v_T) expected regret $\mathcal{R}_T(\{v_i\}_{i=1}^T)$ is bounded above by

$$\left(\gamma + \eta \cdot (e-2)\frac{K+1}{K} \cdot \left(\frac{2K+1}{6} + \frac{\lambda^2}{\gamma}\right) + \frac{\lambda}{K}\right) \cdot T + \frac{\log(K+1)}{\eta}.$$
 (10)

Suppose additionally that $c_1, c_2, c_3 > 0$ are constants. Then there exists a constant c_4 such that, if we set $\gamma = c_1 \cdot (\frac{\log(T)}{T})^{1/3}$, $\eta = c_2 \cdot \gamma^2$, and $K = \lfloor c_3/\gamma \rfloor$, the expected regret $\mathcal{R}_T(\{v_i\}_{i=1}^T)$ is bounded above by

$$c_4 \cdot \log(T)^{1/3} T^{2/3}. \tag{11}$$

COROLLARY 1—Stochastic Upper Bound on Regret of Tempered Exp3 for Social Welfare: Under the assumptions of Theorem 2, suppose additionally that v_i is drawn i.i.d. from some distribution with associated demand function **G**. Then expected regret $\mathcal{R}_T(G)$ is bounded above by the same expressions as in Theorem 2.

The proof of Theorem 2 can again be found in Appendix A.

Tuning. The statement of the theorem leaves the constants c_1, c_2, c_3 in the definition of the tuning parameters unspecified. Suppose we wish to choose the tuning parameters so as to optimize the upper bound obtained in Theorem 2. Ignoring the rounding of K, an approximate solution to this problem is given by

$$\eta = 1/a \cdot \left(\log(T)/T\right)^{2/3},$$

$$\gamma = \lambda \sqrt{(e-2)/a} \cdot \left(\log(T)/T\right)^{1/3},$$

$$K = \sqrt{3\lambda a/(e-2)} \cdot \left(T/\log(T)\right)^{1/3},$$

where

$$a = (9(e-2))^{1/3} (\sqrt{\lambda/3} + \lambda)^{2/3}.$$

This solution is obtained by taking the upper bound in Equation (10), approximating $(K+1)/K \approx 1$ and $(2K+1)/6 \approx K/3$, and solving the first-order conditions with respect to the three tuning parameters. This approximation, and the tuning parameters specified above, yield an approximate upper bound on regret of $6 \cdot \log(T)^{1/3}T^{2/3}$.

Unknown Time Horizon. Note that the proposed tuning depends crucially on knowledge of the time horizon T at which regret is to be evaluated. In order to extend our rate results to the case of unknown time horizons, we can use the so-called doubling trick; cf. Section 2.3 of Cesa-Bianchi and Lugosi (2006): Consider a sequence of epochs (intervals of time periods) of exponentially increasing length, and rerun Algorithm 1 for each time period separately, tuning the parameters over the current epoch length. This construction converts Algorithm 1 into an "anytime algorithm," which enjoys the same regret guarantees of Theorem 2, up to a multiplicative constant factor. Another more efficient strategy to achieve the same goal is to modify Algorithm 1, allowing the parameters η and γ to change at each iteration, and splitting each bin associated with the discretization parameter K whenever more precision is required.

1084



FIGURE 2.—Tempered Exp3 for Social Welfare—numerical example. *Notes*: This figure illustrates the performance of our algorithm for the stochastic case, where v_i is drawn uniformly at random from [0, 1] for all i, the weight λ equals .7, and the tuning parameters are K = 20, $\eta = 0.025$, $\gamma = 0.1$. The left plot shows the cumulative average regret of our algorithm, averaged across 4000 simulations. The right plot shows expected social welfare U(x) as a function of the policy x.

The Extra $\log(T)^{1/3}$ Term. There is a rate discrepancy between our our upper and lower bounds on regret, corresponding to the $\log(T)^{1/3}$ term in our upper bound. We conjecture the existence of an alternative algorithm that can eliminate this extra logarithmic term, albeit at the cost of reduced computational efficiency and a less transparent theoretical analysis. Our conjecture is based on known results for the standard multi-armed bandit problem with K arms. The Exp3 algorithm achieves an upper bound of order $\sqrt{K \log(K)T}$ for this problem, which includes an extra logarithmic factor compared to the known lower bound of order \sqrt{KT} . Exp3 is an instance of the Follow-The-Regularized-Leader (FTRL) algorithm with importance weighting and the negative entropy as the regularizer. It is known that using the $\frac{1}{2}$ -Tsallis entropy as the regularizer in the FTRL algorithm with importance weighting results in regret guarantees of order \sqrt{KT} for the bandit problem (Lattimore and Szepesvári (2020)). However, unlike Exp3, FTRL with Tsallis entropy involves a more complex proof. Analogous statements might be true for our setting.

Numerical Example. For illustration, Figure 2 plots the cumulative average regret of Tempered Exp3 for Social Welfare for the case where v_i is sampled uniformly at random from [0, 1] each time period. Initially, the performance of our algorithm is, by construction, equal to the performance of choosing a policy uniformly at random. Over time, however, the average regret of our algorithm drops by more than half, in this numerical example. Note that the rate at which cumulative regret declines in Figure 2 (for i.i.d. sampling from a fixed distribution) is unrelated to the regret rate of Theorem 2 (for the worst-case sequence of v_i , for each time horizon T).

Alternative Algorithms. Theorem 2 shows that Tempered Exp3 for Social Welfare achieves the lower bound for adversarial regret. The same might be true for other algorithms. Any alternative algorithm that shares this property needs to be randomized. The need for randomization parallels the need for mixed strategies in both static and dynamic zero-sum games; it excludes deterministic algorithms such as UCB. For the bandit setting, the *Tsallis-INF* algorithm (Zimmert and Seldin (2021)), of which Exp3 is a special

case, is furthermore the only algorithm known to be rate optimal in both stochastic and adversarial regimes.

For our adaptive welfare problem, any algorithm that achieves the optimal rate is not only required to randomize; any such algorithm also needs to sample suboptimal policies at a sufficient rate; cf. the proof of Theorem 1. Tempered Exp3 for Social Welfare does so by sampling policies uniformly at random, with probability γ . In the conclusion, we propose a similar modification for Thompson sampling.

A possible improvement to uniform sampling across all policies, as in Tempered Exp3 for Social Welfare, could be to only sample policies uniformly at random from the range of potentially optimal policies: Demand outside this range is irrelevant for welfare comparisons within this range. This idea is implemented in the algorithm that we introduce in Section 4 for the stochastic concave setting.

4. STOCHASTIC REGRET BOUNDS FOR CONCAVE SOCIAL WELFARE

Theorem 1 in Section 3 provides a lower bound proportional to $T^{2/3}$ for adversarial and stochastic regret in social welfare maximization. The proof of this lower bound constructs a distribution for the v_i . This distribution is such that expected social welfare U(x) is nonconcave, as a function of x; two global optima are separated by a region of lower welfare. In order to learn which of two candidates for the globally optimal policy is actually optimal, it is necessary to sample policies in between. These intermediate policies yield lower welfare, and sampling them contributes to cumulative regret. This construction is illustrated in Figure 1.

Given that the construction relies on nonconcavity of expected social welfare, could we achieve lower regret if we knew that social welfare is actually concave? The answer turns out to be yes, for the stochastic setting (in the adversarial setting, cumulative welfare is necessarily nonconcave). One reason is that concavity ensures that the function is unimodal. To estimate the difference in social welfare between two policies, it therefore suffices to sample policies that lie in the interval between them. These in-between policies yield social welfare exceeding the minimum of the two boundary policies. A second reason is that concavity prevents unexpected spikes in social welfare. This property allows us to test carefully chosen triples of points for extended periods, to ensure that one of them is suboptimal, without incurring significant regret.

For the stochastic setting with concave social welfare, we present an algorithm that achieves a bound on regret of order $T^{1/2}$, up to logarithmic terms. Before describing our proposed algorithm, Dyadic Search for Social Welfare, let us formally state the improved regret bounds. The proofs of these lower and upper bounds can be found in the Online Supplement.

THEOREM 3—Lower bound on regret for the concave case: Consider the setup of Section 2. There exists a constant C > 0 such that, for any randomized algorithm for the choice of $x_1, x_2, ...$ and any time horizon $T \in \mathbb{N}$, the following holds.

There exists a distribution μ on [0, 1] with associated demand function G and concave social welfare function U, for which the stochastic cumulative expected regret $\mathcal{R}_T(G)$ is at least $C \cdot T^{1/2}$.

THEOREM 4—Stochastic Upper Bound on Regret of Dyadic Search for Social Welfare: Consider the stochastic setup of Section 2 and time horizon $T \in \mathbb{N}$. If Algorithm 2 is run with confidence parameter $\delta = \frac{1}{T^{5/2}}$, and if the social welfare function U is concave, then the expected regret $\mathcal{R}_T(G)$ is of order at most $T^{1/2}$, up to logarithmic terms. Algorithm 2 Dyadic Search for Social Welfare. **Require:** A confidence parameter $\delta \in (0, 1)$. 1: $I_1 = [0, 1], t_0 = 0, k = 0$ 2: for epochs $\tau = 1, 2, ...$ do Let $c = (\sup I_{\tau} + \inf I_{\tau})/2$, and $d = \sup I_{\tau} - \inf I_{\tau}$. {Subinterval for sampling} 3: 4: if τ is odd then Let $l = c - \frac{1}{4}d$, $r = c + \frac{1}{4}d$. 5: 6: else Let $l = c - \frac{1}{6}d$, $r = c + \frac{1}{6}d$. 7: end if 8: 9: for $t = t_{\tau-1} + 1, t_{\tau-1} + 2, \dots$ do Select $w \in \operatorname{argmax}_{w' \in \{l,c,r,(l,c),(c,r)\}} \Gamma_{t-1}(w')$, {Sampling} 10: breaking ties following the order l, c, r, (l, c), (c, r)if $w \in \{l, c, r\}$ then 11: Set $x_t = w$. 12: 13: else Set $x_t = w_1 + (w_2 - w_1) \cdot \frac{k + 1/2}{n_{t-1}(w_1, w_2) + 1}$, and $k = (k+1) \mod n_{t-1}(w_1, w_2) + 1$. 14: end if 15: Calculate $J_t(l, c)$, $J_t(c, r)$, and $J_t(l, r)$, as in (15) and (16). {Inference} 16: if $\inf(J_t(l, c)) > 0$ or $\inf(J_t(l, r)) > 0$ then 17: let $I_{\tau+1} = I_{\tau} \cap [l, 1]$ and $t_{\tau} = t$ and break {Shrinking the active interval} 18: 19: else if $\sup(J_t(c, r)) \leq 0$ or $\sup(J_t(l, r)) \leq 0$ then 20: let $I_{\tau+1} = I_{\tau} \cap [0, r]$ and $t_{\tau} = t$ and break 21: end if end for 22:

23: end for

Dyadic Search. Our algorithm is based on a modification of dyadic search, as discussed in Bachoc, Cesari, Colomboni, and Paudice (2022a, 2022b). At any point in time, this algorithm maintains an active interval I_{τ} , which contains the optimal policy with high probability. Only policies within this interval are sampled going forward. As evidence accumulates, this interval is trimmed down, by excluding policies that are suboptimal with high probability.

The algorithm proceeds in epochs τ . At the start of each epoch, a subinterval $[l, r] \subset I_{\tau}$ is formed, with mid-point c = (l + r)/2. The points l, c, r are in a dyadic grid, that is, they are of the form $k/2^m$. After sampling from [l, r], we calculate confidence intervals $J_t(l, c), J_t(c, r)$, and $J_t(l, r)$ for the welfare differences $\Delta(l, c), \Delta(c, r)$, and $\Delta(l, r)$, where $\Delta(x, x') = U(x') - U(x)$.

If the confidence interval $J_t(l, c)$ or $J_t(l, r)$ lies above 0, concavity implies that the optimal policy cannot lie to the left of l; we can thus trim the active interval I_τ by dropping all points to the left of l. Symmetrically, if the confidence interval $J_t(c, r)$ or $J_t(l, r)$ lies below 0, we can trim I_τ by dropping all points to the right of r.

Confidence Intervals for Welfare Differences. This procedure requires the construction of confidence intervals for welfare differences of the form

$$\Delta(x, x') = U(x') - U(x) = x' \cdot G(x') - x \cdot G(x) - \lambda \int_x^{x'} G(x'') dx''.$$
(12)

At time t, we estimate demand G(x), for policies x chosen in previous periods, as³

$$\widehat{\boldsymbol{G}}_{t}(x) = \frac{1}{n_{t}(x)} \sum_{i \leq t} y_{i} \cdot \boldsymbol{1}(x_{i} = x), \qquad n_{t}(x) = \sum_{i \leq t} \boldsymbol{1}(x_{i} = x).$$

We similarly estimate integrated demand $\int_x^{x'} G(x'') dx''$ by (x' - x) times the average of realized demand y_i for observations x_i in the open interval (x, x'). We have to be careful, however, to use a sample of x_i that is (approximately) uniformly distributed over this interval. This can be achieved for our dyadic search procedure, as specified in Algorithm 2, by truncating the time index used to estimate this average.⁴ Let

$$s(x, x', t) = \max\left\{s \leq t : \log_2\left(1 + \sum_{i \leq s} \mathbf{1}(x_i \in (x, x'))\right) \in \mathbb{N}\right\}.$$

We define

$$\widehat{G}_{t}(x,x') = \frac{1}{n_{t}(x,x') + 1} \sum_{i \le s(x,x',t)} y_{i} \cdot \mathbf{1}(x_{i} \in (x,x')), \qquad n_{t}(x,x') = \sum_{i \le s(x,x',t)} \mathbf{1}(x_{i} \in (x,x')).$$

At each round, Algorithm 2 maintains estimates for welfare differences among three points l, c, r (for left, center and right, respectively). The estimate of the welfare difference between x' = c and x = l (or between x' = r and x = c) is given by

$$\widehat{\Delta}_{t}(x, x') = x' \cdot \widehat{G}_{t}(x') - x \cdot \widehat{G}_{t}(x) - \lambda \cdot (x' - x) \cdot \widehat{G}_{t}(x, x').$$
(13)

while the estimate of the welfare difference between r and l is given by

$$\widehat{\Delta}_t(l,r) = \widehat{\Delta}_t(l,c) + \widehat{\Delta}_t(c,r).$$
(14)

To construct confidence intervals for $\Delta(x, x')$, we also need to quantify the uncertainty of our demand estimates. We use the following interval half-lengths for confidence intervals for tax revenue at x, and for the private welfare difference between x' and x:

$$\Gamma_t(x) = x \cdot \sqrt{\frac{1}{2n_t(x)} \log\left(\frac{2}{\delta}\right)},$$

$$\Gamma_t(x, x') = \lambda \cdot (x' - x) \cdot \left(\sqrt{\frac{1}{2(n_t(x, x') + 1)} \log\left(\frac{2}{\delta}\right)} + \frac{2}{n_t(x, x') + 1}\right).$$

Using the shorthand $a \pm b = [a - b, a + b]$, our confidence interval for $\Delta(x, x')$, where x' = c and x = l (or x' = r and x = c) is given by

$$J_t(x, x') = \widehat{\Delta}_t(x, x') \pm (\Gamma_t(x') + \Gamma_t(x) + \Gamma_t(x, x')), \qquad (15)$$

1088

³We use the convention 0/0 = 0 and $a/0 = +\infty$ whenever a > 0. Furthermore, every summation over an empty set of indices is understood to have value 0.

⁴The sampling procedure in Algorithm 2 samples sequentially from the dyadic grid in the active interval, refining the grid in subsequent iterations. s(x, x', t) provides a truncation of the time index such that one round of such dyadic sampling has been completed.

while our confidence interval for $\Delta(l, r)$ is given by

$$J_t(l,r) = \widehat{\Delta}_t(l,r) \pm \left(\Gamma_t(r) + \Gamma_t(l) + \Gamma_t(l,c) + \Gamma_t(c,r)\right).$$
(16)

With these preliminaries, we are now ready to state our algorithm, Dyadic Search for Social Welfare, in Algorithm 2.

Before concluding this section, we highlight two features of Algorithm 2. First, two of the three points l, c, r, and the corresponding estimates of demand, are kept from each epoch to the next. Second, estimation of the integral term is performed by querying points following a fixed and balanced design on the dyadic grid—instead of, for example, using a randomized Monte Carlo procedure, which may lead to unbalanced exploration. This implies that the points queried to estimate the integral terms can be easily reused to obtain other integral estimates from each epoch to the next. These two features combined ensure that Algorithm 2 recycles information very efficiently to prune the active interval as quickly as possible, which leads to better regret.

5. INCOME TAXATION

We discuss two extensions of the baseline model of optimal taxation that we introduced in Section 2. These extensions incorporate features that are important in more realistic models of optimal taxation. For both of these extensions, we propose a properly modified version of Algorithm 1. The first extension, discussed in this section, is a variant of the Mirrlees model of optimal income taxation (Mirrlees (1971), Saez (2001, 2002)). The second extension, discussed in Section 6 is a variant of the Ramsey model of commodity taxation (Ramsey (1927)).

Our model of income taxation generalizes our baseline model by allowing for heterogeneous wages w_i , welfare weights $\omega(w_i)$, extensive-margin labor supply responses determined by the cost of participation v_i , and nonlinear income taxes $x_i = \mathbf{x}(w_i)$. Two simplifications are maintained in this model, relative to a more general model of income taxation. First, only extensive margin responses (participation decisions) by individuals are allowed; there are no intensive margin responses (hours adjustments). Second, as in the baseline model of Section 2, there are no income effects. In imposing these assumptions, our model mirrors the model of optimal income taxation discussed in Section II.2 of Saez (2002).

Setup. At each time i = 1, 2, ..., T, one individual arrives who is characterized by (i) a potential wage $w_i \in [0, 1]$, and (ii) an unknown cost of participation $v_i \in [0, 1]$. This individual makes a binary labor supply decision y_i . If they participate in the labor market $(y_i = 1)$, they earn w_i , but pay a tax according to the tax rate $x_i = \mathbf{x}(w_i)$ on their earnings w_i . They furthermore incur a nonmonetary cost of participation v_i .

Their optimal labor supply decision is therefore given by $y_i = \mathbf{1}(v_i \le w_i \cdot (1 - x_i))$, and private welfare equals $\max(w_i \cdot (1 - x_i) - v_i, 0)$. The implied public revenue is equal to the tax on earnings $x_i \cdot w_i$ if $y_i = 1$, and 0 otherwise.

We define social welfare as a weighted sum of public revenue and private welfare, with a weight $\omega(w_i)$ for the latter. Typically, ω is a decreasing function of w, reflecting a preference for redistribution toward those with lower earnings potential; cf. Saez and Stantcheva (2016). Social welfare for time period *i*, as a function of the tax schedule $\mathbf{x}(\cdot)$, is therefore given by

$$U_i(\mathbf{x}(\cdot)) = \underbrace{\mathbf{x}(w_i) \cdot w_i \cdot \mathbf{1}(v_i \leq w_i \cdot (1 - \mathbf{x}(w_i)))}_{\sim}$$

Public revenue

N. CESA-BIANCHI, R. COLOMBONI, AND M. KASY

$$+ \omega(w_i) \cdot \underbrace{\max(w_i \cdot (1 - \mathbf{x}(w_i)) - v_i, 0)}_{\text{Private welfare}}.$$
 (17)

After period *i*, we observe y_i and the tax schedule $\mathbf{x}_i(\cdot)$. If $y_i = 1$, we also observe w_i . Nothing else is observed.⁵

Piecewise Constant Tax Schedules. We next construct a generalization of Algorithm 1 based on piecewise constant tax schedules, with tax rates changing at the grid-points \mathcal{W} , where $0 \in \mathcal{W} \subset [0, 1]$. Formally, define $\lfloor w \rfloor = \max\{w' \in \mathcal{W} : w' \leq w\}$, rounding the wage w down to the nearest grid point in \mathcal{W} ,⁶ Denote $H = |\mathcal{W}|$, and let

$$\mathcal{X}_{\mathcal{W}} = \{ \mathbf{x}(\cdot) : \forall w \in [0, 1], \mathbf{x}(w) = \mathbf{x}(\lfloor w \rfloor) \}.$$

For $w \in \mathcal{W}$ and any $x \in [0, 1]$, denote

$$G_i(w, x) = w_i \cdot \mathbf{1} (v_i \le w_i \cdot (1-x)) \cdot \mathbf{1} (\lfloor w_i \rfloor = w),$$

so that $y_i \cdot w_i = G_i(w_i, \mathbf{x}_i(w_i))$. $G_i(w, x)$ is the individual labor supply function, in monetary units, interacted with an indicator for whether the wage w_i falls into the tax bracket starting at w. With this notation, we can rewrite

$$\max(w_i \cdot (1-x) - v_i, 0) = \int_x^1 G_i(\lfloor w_i \rfloor, x') \, \mathrm{d}x'.$$

For piecewise constant tax rates $\mathbf{x}(\cdot)$, we get

$$U_i(\mathbf{x}(\cdot)) = \sum_{w \in \mathcal{W}} \left[\mathbf{x}(w) \cdot G_i(w, \mathbf{x}(w)) + \omega(w_i) \cdot \int_{\mathbf{x}(w)}^1 G_i(w, x') \, \mathrm{d}x' \right].$$
(18)

Cumulative social welfare is given by $\mathbb{U}_i = \sum_{j \le i} U_i(\mathbf{x}_i(\cdot))$, and we correspondingly define cumulative expected regret, in the adversarial setting, as

$$\mathcal{R}_T = \sup_{\mathbf{x}(\cdot)\in\mathcal{X}_{\mathcal{W}}} E\big[\mathbb{U}_T\big(\mathbf{x}(\cdot)\big) - \mathbb{U}_T\big|\{v_i\}_{i=1}^T, \{w_i\}_{i=1}^T\big].$$

The supremum here is taken over all tax schedules $\mathbf{x}(\cdot)$ that are piecewise constant between the grid points $w \in \mathcal{W}$.

Algorithm. Algorithm 3 generalizes Algorithm 1 to this setting. As before, we form an unbiased estimate \hat{G}_i of G_i using inverse probability weighting, map this estimate into a corresponding estimate \hat{U}_i of U_i , based on Equation (18), and cumulate across time periods to obtain \widehat{U}_i . Note that w_i is observed whenever $y_i = 1$. This implies that the estimate \hat{G}_i is in fact a function of observables, and the same holds for \widehat{U}_i .

1090

⁵It should be noted that in this model we take the transfer x_0 for individuals without other income as given. The effective tax owed by an employed individual equals $\mathbf{x}(w_i) \cdot w_i - x_0$. The "unconditional basic income" x_0 does not affect labor supply, given our assumption that there are no income effects, and it enters social welfare additively. It is therefore without loss of generality to omit x_0 from our model.

⁶Here, we use slightly nonstandard notation, where $\lfloor \cdot \rfloor$ denotes rounding down to the nearest grid-point, rather than the nearest integer.

Algorithm 3 Tempered Exp3 for Optimal Income Taxation.

- **Require:** Tuning parameters K, γ and η , and set of grid points $\mathcal{W} \subset [0, 1]$.
- 1: Calculate evenly spaced grid-points $\mathcal{X} = \{0, \frac{1}{K}, \frac{2}{K}, \dots, 1\}.$
- 2: Initialize $\widehat{\mathbb{G}}_1(w, x) = 0$ and $\widehat{\mathbb{U}}_1(w, x) = 0$ for all $w \in \mathcal{X}$ and all $x \in \mathcal{X}$.
- 3: **for** individual i = 1, 2, ..., T **do**
- 4: For all $x, w \in \mathcal{X}$, set $\lfloor w \rfloor = \max\{w' \in \mathcal{W} : w' \le w\}$, and

{Assignment probabilities}

$$p_i(x|w) = (1 - \gamma) \cdot \frac{\exp(\eta \cdot \widehat{\mathbb{U}}_i(x, \lfloor w \rfloor))}{\sum_{x' \in \mathcal{X}} \exp(\eta \cdot \widehat{\mathbb{U}}_i(x', \lfloor w \rfloor))} + \frac{\gamma}{K+1}.$$
 (19)

5: Draw $A_i \sim U[0, 1]$. For all $w \in [0, 1]$, set

 $\mathbf{x}_{i}(w) = \max\left\{x \in \mathcal{X} : \sum_{x' \in \mathcal{X}, x' < x} p_{i}(x'|w) \le A_{i}\right\},$ (20)

and query y_i accordingly.

6: For all $w \in \mathcal{W}$ and $x \in \mathcal{X}$, set

{Estimated labor supply}

$$\widehat{G}_i(x,w) = y_i \cdot w_i \cdot \frac{\mathbf{1}(\lfloor w_i \rfloor = w, \mathbf{x}_i(w_i) = x)}{p_i(x|w)}.$$
(21)

7: For all
$$w \in \mathcal{W}$$
 and $x \in \mathcal{X}$, set

{Estimated welfare}

$$\widehat{\mathbb{U}}_{i+1}(x,w) = \widehat{\mathbb{U}}_i(x,w) + x \cdot \widehat{G}_i(x,w) + \frac{\omega(w_i)}{K} \cdot \sum_{x' \in \mathcal{X}, x' > x} \widehat{G}_i(x',w).$$
(22)

8: end for

Algorithm 3 keeps track of estimated demand and social welfare for each bin ("tax bracket"), as defined by the grid points $w \in W$. The algorithm then constructs a distribution $p_i(x|w)$ over tax rates $x \in \mathcal{X}$ given w, using the tempered Exp3 distribution. The tax schedule $\mathbf{x}(\cdot)$ is sampled according to these (marginal) distributions of tax rates for each bracket. Though immaterial for the following theorem, we choose the perfectly correlated coupling, across brackets, of these marginal distributions, which is implemented using the random variable A_i in Algorithm 3.

THEOREM 5—Adversarial Upper Bound on Regret of Tempered Exp3 for Optimal Income Taxation: Consider the setup of Section 5, and Algorithm 3. Assume that $(K+1)\eta < \gamma$, and that $\omega(w) \le 1$ for all w.

Then for any sequence (v_1, \ldots, v_T) expected regret $\mathcal{R}_T(\{v_i\}_{i=1}^T)$ is bounded above by

$$\left(\gamma + \eta \cdot (e-2)\frac{K+1}{K} \cdot \left(\frac{2K+1}{6} + \frac{1}{\gamma}\right) + \frac{1}{K}\right) \cdot T + \frac{H\log(K+1)}{\eta}.$$
 (23)

Suppose additionally⁷ that $K = c_1 \cdot (T/H)^{1/3}$, $\gamma = c_2/(K+1)$, and $\eta = c_3/(K+1)^2$, for some constants c_1, c_2, c_3 . Then expected regret $\mathcal{R}_T(\{v_i\}_{i=1}^T)$ is bounded above by

$$c_4 \cdot H^{1/3} \cdot \log(T)^{1/3} T^{2/3},$$
 (24)

for some constant c_4 .

6. COMMODITY TAXATION

In this section, we generalize our baseline model of optimal taxation to a model of commodity taxation with multiple goods $j \in \{1, ..., k\}$ and continuous demand functions $y_i(x) \in [0, 1]^k$, where $x \in [0, 1]^k$ is a vector of tax rates. We again assume that there are no income effects. Our setup is a version of the classic Ramsey model (Ramsey (1927)). We propose a generalization of Tempered Exp3 for Social Welfare to this setting. In the following, we use $\langle x, y \rangle$ to denote the Euclidean inner product between x and y.

Setup. At each time i = 1, 2, ..., T, one individual arrives who is characterized by a utility function $u_i : [0, 1]^k \to \mathcal{R}$. This individual is exposed to a tax vector $x_i \in [0, 1]^k$, and makes a continuous consumption decision y_i . Public revenue is given by $\langle x_i, y_i \rangle$. Private utility is given by $u_i(y_i)$ plus their consumption of a numeraire good, which has price normalized to 1 and enters utility additively. The individual consumption choice y_i costs $\langle x_i + p, y \rangle$, where p is the (exogenously given) vector of pre-tax prices. This cost of purchasing y_i reduces the consumption of the numeraire good. The optimal individual decision is therefore given by

$$y_i = G_i(x_i) = \underset{y \in [0,1]^k}{\operatorname{argmax}} [u_i(y) - \langle x_i + p, y \rangle].$$
(25)

The implied private welfare is

$$v_i(x) = v_0 + \max_{y \in [0,1]^k} \left[u_i(y) - \langle x + p, y \rangle \right],$$

where we have added a constant v_0 , chosen such that $v_i(0) = 0$; this is just a normalization to simplify notation below.

We define social welfare as a weighted sum of public revenue and private welfare, with a weight λ for the latter. Social welfare for time period *i*, as a function of the tax vector *x*, is therefore given by

$$U_i(x_i) = \underbrace{\langle x_i, y_i \rangle}_{\text{Public revenue}} + \lambda \cdot \underbrace{v_i(x_i)}_{\text{Private welfare}} .$$
(26)

After period *i*, we observe y_i and the tax vector x_i . Nothing else is observed. Algorithm 4 adapts our approach to this setting. This algorithm requires a mapping Π from (estimated) demand to welfare.

Mapping Demand to Welfare. By the envelope theorem (Milgrom and Segal (2002)),

$$\nabla_x v_i(x) = G_i(x).$$

⁷ for simplicity, we assume that in the following tuning K is an integer. If not, round K to the closest integer.

Algorithm 4 Tempered Exp3 for Commodity Taxation.

Require: Tuning parameters K, γ and η .

- 1: Calculate the set of evenly spaced grid-points $\mathcal{X} = \{0, \frac{1}{\kappa}, \dots, 1\}^k$ and initialize $\widehat{\mathbb{G}}_1(x) = 0$ for all grid points.
- 2: **for** individual i = 1, 2, ..., T **do**
- For all $x \in \mathcal{X}$, set 3:

$$\widehat{\mathbb{U}}_{i}(x) = \langle x_{i}, \widehat{\mathbb{G}}_{i} \rangle + \lambda \cdot \widehat{v}_{i}(x_{i}).$$
(28)

For all $x \in \mathcal{X}$, set 4:

{Estimated welfare}

$$p_i = (1 - \gamma) \cdot \frac{\exp(\eta \cdot \mathbb{U}_i(x))}{\sum_{x'} \exp(\eta \cdot \widehat{\mathbb{U}}_i(x'))} + \frac{\gamma}{(K+1)^k}.$$
(29)

- Choose x_i at random according to the probability distribution p_i , and query y_i ac-5: cordingly. {Estimated demand}
- For all $x \in [0, 1]^k$, set 6:

$$\widetilde{\mathbb{G}}_{i+1}(x) = \widehat{\mathbb{G}}_{i}(x) + y_{i} \cdot \frac{\mathbf{1}(x_{i} = \lfloor x \rfloor)}{p_{i}}$$
$$\widehat{v}_{i+1}(x) = \Pi(\widetilde{\mathbb{G}}_{i+1})$$
$$\widehat{\mathbb{G}}_{i+1}(x) = \nabla_{x}\widehat{v}_{i+1}(x).$$
(30)

7: end for

Let \mathcal{V} be the set of differentiable functions v on $[0, 1]^k$ such that $\nabla_x v \in L^2$, and such that v(0) = 0. Consider the following operator, mapping the demand function G into the corresponding indirect utility function v:

$$\Pi(G(\cdot)) \in \underset{v(\cdot)\in\mathcal{V}}{\operatorname{argmin}} \int_{[0,1]^k} \left\| \nabla_x v(x) - G(x) \right\|^2 \mathrm{d}x$$
(27)

We can think of the operator Π as combining two operators. First, the function G is projected on the subspace of functions on $[0, 1]^k$, which can be written as the gradient of some function v. Second, the projected G is then integrated to get v(x) for any x. Integration here is along some curve in $[0, 1]^k$ from 0 to x. Given the first projection, the choice of curve does not matter for the resulting function v. A formal analysis of Tempered Exp3 for Commodity Taxation would need to prove existence of the projection. We leave such a formal analysis, including lower and upper regret bounds, for future research.

7. CONCLUSION

Possible Applications. The setup introduced in Section 2 was deliberately stylized, to allow for a clear exposition of the conceptual issues that arise when adaptively maximizing social welfare. The algorithm that we proposed for this setup, and the generalizations

discussed later in the paper, are nonetheless directly practically relevant. They remain appropriate in economic settings that are considerably more general than the setting described by our model.

The reasons for this generality have been elucidated by the public finance literature, cf. Chetty (2009), building on the generality of the envelope theorem; cf. Milgrom and Segal (2002), Sinander (2022). By the envelope theorem, the welfare impact of a marginal tax change on private welfare can be calculated ignoring any behavioral responses to the tax change. This holds in generalizations of our setup that allow for almost arbitrary action spaces (including discrete and continuous, multidimensional, and dynamic actions), and for arbitrary preference heterogeneity. The expressions for social welfare that justify our algorithms remain unchanged under such generalizations. That said, the validity of these expressions does require the absence of income effects and of externalities. If there are income effects or externalities, the algorithms need to be modified.

Our approach is motivated by applications of algorithmic decision-making for public policy, where a policymaker cares about welfare, but also faces a government budget constraint. Possible application domains of our algorithm include the following. In public health and development economics, field experiments such as Cohen and Dupas (2010) vary the level of a subsidy for goods such as insecticide-treated bed nets (ITNs), estimating the impact on demand. Our algorithm could be used to find the optimal subsidy level quickly and apply it to experimental participants. A term capturing positive externalities of the use of ITNs could be added to social welfare, leaving the algorithm otherwise unchanged. In educational economics, many studies evaluate the impact of financial aid on college enrollment (Dynarski, Page, and Scott-Clayton (2023)). An adaptive experiment might vary the level of aid provided, where aid is conditional on college attendance and conditional on pre-determined criteria of need or merit. In such an experiment, a variant of our algorithm for optimal income taxation might be used, where the welfare weights ω are a function of need or merit, and the outcome y is college attendance. In *environmental* economics, many experiments (e.g., Lee, Miguel, and Wolfram (2020)) study the impact of electricity pricing on household electricity consumption. Once again, our baseline algorithm (for binary household decisions about connecting to the grid) or our algorithm for commodity taxation (for continuous household decisions about consumption levels) could be applied, to learn optimal prices, taking into account both distributional considerations and externalities.

These examples are all drawn from public policy, where there is an intrinsic concern for social welfare. This contrasts with commercial applications, where the goal is typically to maximize (directly observable) profits by monopolist pricing (den Boer (2015)), or more generally by reserve price setting in auctions (Nedelec, Calauzènes, El Karoui, and Perchet (2022). Adaptive pricing algorithms are used in applications such as online ad auctions. A concern for welfare might enter in such commercial settings if there is a participation constraint that needs to be satisfied for consumers. Suppose, for example, that consumers or service providers need to first sign up for a platform, say for e-commerce or for gig work, and then repeatedly engage in transactions on this platform. To sign up in the first place, their expected welfare needs to exceed their outside option. This constraint might then enter the platform provider's objective, in Lagrangian form, adding a term for welfare, and leading to objectives such as those maximized by our algorithms.

An Alternative Approach: Thompson Sampling. The main algorithm proposed in this paper, Tempered Exp3 for Social Welfare, is designed to perform well in the adversarial setting. In the construction of this algorithm, no probabilistic assumptions were made

1095

about the distribution of v_i . In the stochastic framework, a sampling distribution is assumed, for instance, that the v_i be i.i.d. over time. The Bayesian framework completes this by assuming a prior distribution over the parameters, which govern the sampling distribution.

One popular heuristic for adaptive policy choice in the Bayesian framework is Thompson sampling (Thompson (1933), Russo, Van Roy, Kazerouni, Osband, and Wen (2018)), also known as probability matching, which assigns a policy with probability equal to the posterior probability that this policy is optimal. In our setting, Thompson sampling could be implemented as follows. First, form a posterior for the demand function G(x) = E[y|x], based on all the data available from previous periods j < i. Sample one draw $\tilde{G}(\cdot)$ from this posterior. Map this draw into a draw $\tilde{U}(\cdot)$ from the posterior for $U(\cdot)$ via $\tilde{U}(x) = x \cdot \tilde{G}(x) + \lambda \cdot \int_x^1 \tilde{G}(x') dx'$. Find the maximizer $x_i = \operatorname{argmax}_x \tilde{U}(x)$. This is the policy recommended by Thompson sampling. We conjecture that this algorithm will outperform random assignment, but will underexplore relative to the optimal algorithm. Adding further forced exploration to this algorithm might improve cumulative welfare. A formal analysis of algorithms of this type is left for future research.

A natural class of priors for G are Gaussian process priors (Williams and Rasmussen (2006)). If outcomes y are conditionally normal (rather than binary, as in our baseline model), then the posterior for demand is available in closed form, and the posterior mean is equal to the best linear predictor given past outcomes y_i . Furthermore, since social welfare is a linear transformation of demand, the posterior for U is then also linear and available in closed form. For details, see Kasy (2018).

APPENDIX A: PROOFS

A.1. Theorem 1 (Lower Bound on Regret)

Defining a Family of Distributions for v. Recall that, for each $\epsilon \in [-1, 1]$, the probability distribution μ^{ϵ} is defined as the probability measure supported on (1/4, 1/2, 3/4, 1) with masses $(a, (1 + \epsilon) \cdot b, (1 - \epsilon) \cdot b, 1 - a - 2 \cdot b)$, where

$$a := \frac{(1-\lambda) \cdot (136-99 \cdot \lambda)}{2 \cdot (4-3 \cdot \lambda) \cdot (24-17 \cdot \lambda)}, \qquad b := \frac{1-\lambda}{2 \cdot (24-17 \cdot \lambda)}.$$

Furthermore, for each $\epsilon \in [-1, 1]$, recall that G^{ϵ} and U^{ϵ} are respectively the demand function and the expected social welfare associated to μ^{ϵ} (see Figure 1 for an illustration). Let $v_1, v_2, \dots \in [0, 1]$ be the sequence of individual valuations. For each $\epsilon \in [-1, 1]$, consider a distribution P^{ϵ} such that the individual valuations v_1, v_2, \dots form a P^{ϵ} -i.i.d. sequence (independent of the randomization used by the algorithm) with common distribution μ^{ϵ} .

Explicit Lower Bound on Regret That Will Be Proven. Define

$$c_1 := \frac{\lambda}{4} \cdot b, \qquad c_2 := \frac{1}{8} \cdot \frac{1 - \lambda}{4 - 3 \cdot \lambda}, \qquad c_3 := b \cdot \sqrt{\frac{2}{a \cdot (1 - a - 2 \cdot b)}}.$$

We will prove that, for any randomized algorithm and any time horizon $T \in \mathbb{N}$, there exists $\epsilon \in [-1, 1]$ such that

$$\mathcal{R}_T(\boldsymbol{G}^{\boldsymbol{\epsilon}}) \geq C \cdot T^{2/3},$$

where

$$C := \min\left(\frac{c_1^2 \cdot c_3^2}{c_2}, \frac{c_2}{2}, \frac{1}{16} \cdot \sqrt[3]{\frac{c_1^2 \cdot c_2}{c_3^2}}\right)$$

= $\min\left(\frac{\lambda^2 \cdot (4 - 3 \cdot \lambda)^3}{8 \cdot (136 - 99 \cdot \lambda) \cdot (26 - 19 \cdot \lambda)}, \frac{\lambda^{2/3} \cdot (1 - \lambda)^{4/3} \cdot (136 - 99 \cdot \lambda)^{1/3} \cdot (26 - 19 \cdot \lambda)^{1/3}}{128 \cdot (4 - 3 \cdot \lambda) \cdot (24 - 17 \cdot \lambda)^{4/3}}\right) > 0.$ (31)

Fix a randomized algorithm to choose the policies $x_1, x_2, ...,$ and fix a time horizon $T \in \mathbb{N}$.

Number of Mistakes and Lower Bound on Regret. We need to count the random number of times the algorithm has played in the regions (1/2, 3/4], [0, 1/2] and (3/4, 1] up to time *T*. This can be done relying on the following random variables:

$$n_1 := \sum_{i=1}^T \mathbf{1}_{(1/2,3/4]}(x_i), \qquad n_2 := \sum_{i=1}^T \mathbf{1}_{[0,1/2]}(x_i), \qquad n_3 := \sum_{i=1}^T \mathbf{1}_{(3/4,1]}(x_i).$$

Notice that since the intervals (1/2, 3/4], [0, 1/2] and (3/4, 1] form a partition of [0, 1], we have that

$$n_1 + n_2 + n_3 = T \tag{32}$$

For each $\epsilon \in [-1, 1]$, denote by E^{ϵ} the expectation taken with respect to the distribution P^{ϵ} . Notice that, for each $\epsilon \in [-1, 1]$, the expected regret when the underlying distribution is P^{ϵ} equals

$$\mathcal{R}_{T}(\boldsymbol{G}^{\boldsymbol{\epsilon}}) = T \cdot \sup_{\boldsymbol{x} \in [0,1]} \boldsymbol{U}^{\boldsymbol{\epsilon}}(\boldsymbol{x}) - \sum_{i=1}^{T} E^{\boldsymbol{\epsilon}}(\boldsymbol{U}^{\boldsymbol{\epsilon}}(\boldsymbol{x}_{i})).$$
(33)

Algebraic calculations show that, for each $\epsilon \in [-1, 1]$,

$$\max_{x \in (1/2, 3/4]} U^{\epsilon}(x) = U^{\epsilon}(3/4), \qquad \max_{x \in [0, 1/2]} U^{\epsilon}(x) = U^{\epsilon}(1/4),$$
(34)

$$\max_{x \in (3/4,1]} U^{\epsilon}(x) = U^{\epsilon}(1), \qquad (C^{\epsilon})$$

and
$$U^{\epsilon}(1) - U^{\epsilon}(1/4) = c_1 \cdot \epsilon.$$
 (35)

Further calculations show also that

$$\min_{\epsilon \in [-1,1]} \min(U^{\epsilon}(1/4), U^{\epsilon}(1)) = U^{1}(1/4), \qquad \max_{\epsilon \in [-1,1]} \max_{x \in (1/2,3/4]} U^{\epsilon}(x) = U^{-1}(3/4), \quad (36)$$

and
$$U^{1}(1/4) - U^{-1}(3/4) = c_2.$$
 (37)

Equations (34), (35), (36), and (37) imply that

$$\sup_{x \in [0,1]} U^{\epsilon}(x) = U^{\epsilon}(1) \quad \text{if } \epsilon \in [0,1].$$
(38)

1096

It follows that, if $\epsilon \in [0, 1]$,

$$\mathcal{R}_{T}(\boldsymbol{G}^{\epsilon}) \stackrel{(33)}{=} T \cdot \sup_{x \in [0,1]} \boldsymbol{U}^{\epsilon}(x) - \sum_{i=1}^{T} E^{\epsilon}(\boldsymbol{U}^{\epsilon}(x_{i}))$$

$$\stackrel{(38)}{=} T \cdot \boldsymbol{U}^{\epsilon}(1) - \sum_{i=1}^{T} E^{\epsilon}(\boldsymbol{U}^{\epsilon}(x_{i}) \cdot (\mathbf{1}_{(1/2,3/4]}(x_{i}) + \mathbf{1}_{[0,1/2]}(x_{i}) + \mathbf{1}_{(3/4,1]}(x_{i}))))$$

$$\stackrel{(34)}{\geq} T \cdot \boldsymbol{U}^{\epsilon}(1) - \sum_{i=1}^{T} E^{\epsilon}(\boldsymbol{U}^{\epsilon}(^{3/4}) \cdot \mathbf{1}_{(1/2,3/4]}(x_{i}))$$

$$+ \boldsymbol{U}^{\epsilon}(^{1/2}) \cdot \mathbf{1}_{[0,1/2]}(x_{i}) + \boldsymbol{U}^{\epsilon}(1) \cdot \mathbf{1}_{(3/4,1]}(x_{i}))$$

$$\stackrel{(32)}{=} (\boldsymbol{U}^{\epsilon}(1) - \boldsymbol{U}^{\epsilon}(^{3/4})) \cdot E^{\epsilon}(n_{1}) + (\boldsymbol{U}^{\epsilon}(1) - \boldsymbol{U}^{\epsilon}(^{1/4})) \cdot E^{\epsilon}(n_{2})$$

$$\stackrel{(36)}{\geq} (\boldsymbol{U}^{1}(^{1/4}) - \boldsymbol{U}^{-1}(^{3/4})) \cdot E^{\epsilon}(n_{1}) + (\boldsymbol{U}^{\epsilon}(1) - \boldsymbol{U}^{\epsilon}(^{1/4})) \cdot E^{\epsilon}(n_{2})$$

$$\stackrel{(37)}{=} c_{2} \cdot E^{\epsilon}(n_{1}) + (\boldsymbol{U}^{\epsilon}(1) - \boldsymbol{U}^{\epsilon}(^{1/4})) \cdot E^{\epsilon}(n_{2})$$

$$\stackrel{(35)}{=} c_{2} \cdot E^{\epsilon}(n_{1}) + c_{1} \cdot \epsilon \cdot E^{\epsilon}(n_{2})$$

$$(39)$$

Notice that inequality (39) quantifies how much regret the algorithm is going to suffer in terms of the expected number of times it plays in the wrong regions, when the demand function is G^{ϵ} and $\epsilon > 0$.

In the same way inequality (39) was proven, we can prove that, if $\epsilon \in [0, 1]$,

$$\mathcal{R}_{T}(\boldsymbol{G}^{-\epsilon}) \geq c_{2} \cdot E^{-\epsilon}(n_{1}) + c_{1} \cdot \boldsymbol{\epsilon} \cdot E^{-\epsilon}(n_{3}) \geq c_{1} \cdot \boldsymbol{\epsilon} \cdot E^{-\epsilon}(n_{3}), \tag{40}$$

which again quantifies how much regret the algorithm is going to suffer in terms of the expected number of times it plays in the wrong regions, when the demand function is $G^{-\epsilon}$ and $\epsilon > 0$.

Intuition for the Remainder of the Proof. At high level, inequalities (39) and (40) tell us that, if $|\epsilon|$ is not negligible, the algorithm has to play a substantially different number of times in the region ($^{3}/_{4}$, 1], depending on the sign of ϵ , not to suffer significant regret when the demand function is G^{ϵ} . The crucial idea is that the only way for the algorithm to present this different behavior is by playing in the only informative region about the sign of ϵ , that is, the region ($^{1}/_{2}$, $^{3}/_{4}$]. However, as shown in (39), selecting policies in this region comes at a cost in terms of regret. To relate quantitatively the number of times the algorithm has to play in this costly region with the difference in the expected number of times the algorithm selects policies in the region ($^{3}/_{4}$, 1] is the last missing ingredient that we can obtain relying on information theoretic techniques: It can be proved (and a formal proof is provided in the Online Supplement, in Section B.1, that, for each $\epsilon \in [0, 1]$,

$$E^{-\epsilon}(n_3) \ge E^{\epsilon}(n_3) - c_3 \cdot \epsilon \cdot T \cdot \sqrt{E^{\epsilon}(n_1)}.$$
(41)

Now, if the algorithm is going to suffer low regret when $\epsilon > 0$, then by (39) we have an upper bound on the number of times the algorithm plays in the region (1/2, 3/4] and a lower bound on the number of times it plays in the region (3/4, 1], whenever $\epsilon > 0$. In

turn, by (41), this gives a lower bound on the number of times the algorithm plays in the suboptimal region ($^{3}/_{4}$, 1] when $\epsilon < 0$. Then, relying on (40), we have an explicit lower bound on how much regret the algorithm is going to suffer when $\epsilon < 0$. We will now carry out this plan—and prove the theorem—as follows.

Low Regret Cannot Be Achieved for Both Positive and Negative ϵ . To get a contradiction, suppose that

$$\forall \epsilon \in [-1, 1] \quad \mathcal{R}_T(G^{\epsilon}) < C \cdot T^{2/3}.$$
(42)

It follows from (39) that, for each $\epsilon \in [0, 1]$,

$$E^{\epsilon}(n_1) \stackrel{(39)}{\leq} \frac{\mathcal{R}_T(\boldsymbol{G}^{\epsilon})}{c_2} \stackrel{(42)}{\leq} \frac{C}{c_2} \cdot T^{2/3}, \qquad E^{\epsilon}(n_2) \stackrel{(39)}{\leq} \frac{\mathcal{R}_T(\boldsymbol{G}^{\epsilon})}{c_1 \cdot \epsilon} \stackrel{(42)}{\leq} \frac{C}{c_1 \cdot \epsilon} \cdot T^{2/3}.$$
(43)

This implies, relying also on (40) and (41), that for each $\epsilon \in [0, 1]$ we have

$$\mathcal{R}_{T}(\boldsymbol{G}^{-\epsilon}) \stackrel{(40)}{\geq} c_{1} \cdot \boldsymbol{\epsilon} \cdot E^{-\epsilon}(n_{3}) \stackrel{(41)}{\geq} c_{1} \cdot \boldsymbol{\epsilon} \cdot \left(E^{\epsilon}(n_{3}) - c_{3} \cdot \boldsymbol{\epsilon} \cdot T \cdot \sqrt{E^{\epsilon}(n_{1})}\right)$$

$$\stackrel{(32)}{=} c_{1} \cdot \boldsymbol{\epsilon} \cdot \left(T - E^{\epsilon}(n_{1}) - E^{\epsilon}(n_{2}) - c_{3} \cdot \boldsymbol{\epsilon} \cdot T \cdot \sqrt{E^{\epsilon}(n_{1})}\right)$$

$$\stackrel{(43)}{\geq} c_{1} \cdot \boldsymbol{\epsilon} \cdot \left(T - \frac{C}{c_{2}} \cdot T^{2/3} - \frac{C}{c_{1} \cdot \boldsymbol{\epsilon}} \cdot T^{2/3} - c_{3} \cdot \boldsymbol{\epsilon} \cdot T \cdot \sqrt{\frac{C}{c_{2}}} \cdot T^{2/3}\right)$$

$$= c_{1} \cdot \boldsymbol{\epsilon} \cdot \left(1 - \frac{C}{c_{2}} \cdot T^{-1/3} - \frac{C}{c_{1} \cdot \boldsymbol{\epsilon}} \cdot T^{-1/3} - c_{3} \cdot \boldsymbol{\epsilon} \cdot T^{1/3} \cdot \sqrt{\frac{C}{c_{2}}}\right) \cdot T. \quad (44)$$

Pick $\epsilon := T^{-1/3} \cdot \sqrt{\frac{\sqrt{C \cdot c_2}}{c_1 \cdot c_3}}$. First, note that since $0 < C \leq \frac{c_1^{-2} \cdot c_3^2}{c_2}$ we have that $\epsilon \in (0, 1]$. Substituting this value of ϵ in (44) leads to

$$C \cdot T^{2/3} \stackrel{(42)}{\geq} \mathcal{R}_{T} (\mathbf{G}^{-\epsilon})$$

$$\stackrel{(44)}{\geq} \sqrt{\frac{\sqrt{C \cdot c_{2} \cdot c_{1}}}{c_{3}}} \cdot \left(1 - \frac{C}{c_{2}} \cdot T^{-1/3} - 2 \cdot \sqrt{\frac{c_{3}}{c_{1} \cdot \sqrt{c_{2}}}} \cdot C^{3/4}\right) \cdot T^{2/3}$$

$$\stackrel{(31)}{\geq} \frac{1}{2} \cdot \sqrt{\frac{\sqrt{C \cdot c_{2} \cdot c_{1}}}{c_{3}}} \cdot \left(1 - 4 \cdot \sqrt{\frac{c_{3}}{c_{1} \cdot \sqrt{c_{2}}}} \cdot C^{3/4}\right) \cdot T^{2/3}$$

$$\stackrel{(31)}{\geq} \frac{1}{4} \cdot \sqrt{\frac{\sqrt{C \cdot c_{2} \cdot c_{1}}}{c_{3}}} \cdot T^{2/3}, \qquad (45)$$

where the second to last inequality follows from $C \le \frac{c_2}{2}$, while the last inequality follows from $C \le \frac{1}{16} \sqrt[3]{\frac{c_1^2 \cdot c_2}{c_3^2}}$. Rearranging inequality (45) leads to the contradiction

$$C^{(45)} \left(\frac{1}{4} \cdot \sqrt{\frac{c_1 \cdot \sqrt{c_2}}{c_3}}\right)^{4/3} = \frac{1}{8} \cdot \sqrt[3]{\frac{2 \cdot c_1^2 \cdot c_2}{c_3^2}} > \frac{1}{16} \cdot \sqrt[3]{\frac{c_1^2 \cdot c_2}{c_3^2}} \stackrel{(31)}{\ge} C.$$

1098

Since (42) leads to a contradiction, it follows that there exists $\epsilon \in [-1, 1]$ such that $\mathcal{R}_T(\mathbf{G}^{\epsilon}) \geq C \cdot T^{2/3}$. Given that the time horizon T and the randomized algorithm were arbitrarily fixed, the theorem is proved.

A.2. Theorem 2 (Adversarial Upper Bound on Regret)

The proof of this theorem builds upon the proof of Theorem 6.5 in Cesa-Bianchi and Lugosi (2006). Relative to this theorem, we need to additionally consider the discretization error introduced by Algorithm 1, and explicitly control the variance of estimated welfare.

Recall our notation \mathbb{U} and $\mathbb{U}(x)$ for realized cumulative welfare, and for cumulative welfare for the counterfactual, fixed policy x. We further abbreviate $\mathbb{U}_{Tk} = \mathbb{U}(\tilde{x}_k)$. Throughout this proof, the sequence $\{v_i\}_{i=1}^T$ is given and conditioned on in any expectations.

1. Discretization

Recall that $U_i(x) = x \cdot \mathbf{1}(x \le v_i) + \lambda \cdot \max(v_i - x, 0)$. Let

$$\tilde{v}_i = \max_k \{ \tilde{x}_k : \tilde{x}_k \le v_i \}$$

(this is v_i rounded down to the next grid point \tilde{x}_k), and denote

$$ilde{U}_i(x) = x \cdot \mathbf{1}(x \le v_i) + \lambda \cdot \max(\tilde{v}_i - x, 0),$$

 $ilde{\mathbb{U}}_i(x) = \sum_{j \le i} \tilde{U}_j(x),$

as well as $\tilde{\mathbb{U}}_{ik} = \tilde{\mathbb{U}}_i(\tilde{x}_k)$. Then it is immediate that $\tilde{U}_i(x) \le U_i(x)$,

$$\sup_{x} \left| \tilde{U}_{i}(x) - U_{i}(x) \right| \leq \frac{\lambda}{K},$$

and $\operatorname{argmax}_{x} \tilde{\mathbb{U}}_{i}(x) \in \{\tilde{x}_{1}, \ldots, x_{K+1}\}$ and, therefore,

$$\max_{k} \tilde{\mathbb{U}}_{ik} \geq \sup_{x} \mathbb{U}_{i}(x) - i \cdot \frac{\lambda}{K}$$

2. Unbiasedness

At the end of period *i*, \widehat{G}_k is an unbiased estimator of $\sum_{j \le i} \mathbf{1}(\widetilde{x}_k \le v_j)$ for all *k*. Therefore, $E[\widehat{U}_{ik}] = \widetilde{U}_{ik}$ for all *i* and *k*.

3. Upper bound on optimal welfare Define $W_i = \sum_k \exp(\eta \cdot \widehat{\mathbb{U}}_{ik})$, and $q_{ik} = \exp(\eta \cdot \widehat{\mathbb{U}}_{ik})/W_i$. It is immediate that

$$E[\log W_T] \ge \eta \cdot E\left[\max_k \widehat{\mathbb{U}}_{Tk}\right] \ge \eta \cdot \max_k E[\widehat{\mathbb{U}}_{Tk}] = \eta \cdot \max_k \widetilde{\mathbb{U}}_{Tk}.$$

Furthermore,

$$E[\log W_T] = \sum_{0 \le i < T} E\left[\log\left(\frac{W_{i+1}}{W_i}\right)\right] + \log(W_0).$$

Given our initialization of the algorithm, $log(W_0) = log(K + 1)$.

4. Lower bound on estimated welfare

Denote $\widehat{U}_{ik} = \widetilde{x}_k \cdot \widehat{H}_k + \frac{\lambda}{K} \cdot \sum_{k'>k} \widehat{H}_{k'}$, where $\widehat{H}_k = \frac{y_i}{p_{ik}} \cdot \mathbf{1}(k_i = k)$, so that $\widehat{\mathbb{U}}_{ik} = \sum_{j < i} \widehat{U}_{jk}$, and $E[\widehat{U}_{jk}] = U_i(\widetilde{x}_k)$. By definition of W_i ,

$$\log\left(\frac{W_{i+1}}{W_i}\right) = \log\left(\sum_k q_{ik} \cdot \exp(\eta \cdot \widehat{U}_{ik})\right).$$

Since $p_k \ge \gamma/(K+1)$ for all k, $\widehat{U}_{ik} \in [0, (K+1)/\gamma]$ for all i and k and, therefore, $\eta \cdot \widehat{U}_{ik} \le (K+1) \cdot \eta/\gamma \le 1$ (where the last inequality holds by assumption). Using $\exp(a) \le 1 + a + (e-2)a^2$ for any $a \le 1$ yields

$$\exp(\eta \widehat{U}_{ik}) \leq 1 + \eta \cdot \widehat{U}_{ik} + (e-2) \cdot (\eta \cdot \widehat{U}_{ik})^2$$

Therefore,

$$\log\left(\frac{W_{i+1}}{W_i}\right) \leq \log\left(\sum_{k} q_{ik} \cdot \left(1 + \eta \cdot \widehat{U}_{ik} + (e-2) \cdot (\eta \cdot \widehat{U}_{ik})^2\right)\right)$$
$$\leq \eta \cdot \sum_{k} q_{ik} \cdot \widehat{U}_{ik} + (e-2) \cdot \eta^2 \cdot \sum_{k} q_{ik} \cdot \widehat{U}_{ik}^2$$

The second inequality follows from $\log(1 + x) \le x$.

5. Connecting the first-order term to welfare

Note that, by definition, $q_{ik} = (p_{ik} - \frac{\gamma}{\kappa+1})/(1-\gamma)$. Therefore,

$$\sum_{k} q_{ik} \cdot \widehat{U}_{ik} = \frac{1}{1-\gamma} \sum_{k} p_{ik} \cdot \widehat{U}_{ik} - \frac{\gamma}{(1-\gamma)(K+1)} \cdot \sum_{k} \widehat{U}_{ik},$$

and thus

$$E\left[\sum_{k} q_{ik} \cdot \widehat{U}_{ik}\right] \leq \frac{1}{1-\gamma} E\left[\widetilde{U}_{i}(x_{i})\right],$$

where we have used the fact that $0 \le \tilde{U}_k \le 1$ for all k, given our definition of \tilde{U} , and the fact that k_i is distributed according to p_{ik} , by construction.

6. Bounding the second moment of estimated welfare It remains to bound the term $E[\sum_{k} q_{ik} \cdot \hat{U}_{ik}^{2}]$. As in the preceding item, we have

$$\sum_{k} q_{ik} \cdot \widehat{U}_{ik}^2 \leq \frac{1}{1-\gamma} \sum_{k} p_{ik} \cdot \widehat{U}_{ik}^2.$$

We can rewrite

$$\widehat{U}_{ik} = \left(\widetilde{x}_k \cdot \mathbf{1}(k_i = k) + \frac{\lambda}{K} \cdot \mathbf{1}(k_i > k)\right) \cdot \frac{y_i}{p_{ik_i}}.$$

Bounding $y_i \le 1$ immediately gives

$$E_i[\widehat{U}_{ik}^2] \le \frac{\widetilde{x}_k^2}{p_{ik}} + \left(\frac{\lambda}{K}\right)^2 \cdot \sum_{k' > k} \frac{1}{p_{ik'}}$$

and, therefore,

$$E_{i}\left[\sum_{k} p_{ik} \cdot \widehat{U}_{ik}^{2}\right] \leq \sum_{k} \widetilde{x}_{k}^{2} + \left(\frac{\lambda}{K}\right)^{2} \cdot \sum_{k} \sum_{k' > k} \frac{p_{ik}}{p_{ik'}}$$
$$\leq \sum_{k} \left(\frac{k}{K}\right)^{2} + \left(\frac{\lambda}{K}\right)^{2} \cdot \sum_{k} p_{ik} \sum_{k' \neq k} \frac{K+1}{\gamma}$$
$$= \frac{K(K+1)(2K+1)}{6K^{2}} + \frac{\lambda^{2}}{\gamma} \frac{K+1}{K}$$
$$= \frac{K+1}{K} \cdot \left(\frac{2K+1}{6} + \frac{\lambda^{2}}{\gamma}\right).$$

7. Collecting inequalities

Combining the preceding items, we get

$$\begin{split} \eta \cdot \left(\sup_{x} \mathbb{U}(x) - T \cdot \frac{\lambda}{K} \right) \\ &\leq \eta \cdot \max_{k} \tilde{\mathbb{U}}_{Tk} \leq E[\log W_{T}] \quad (\text{Item 1}) \\ &= \sum_{0 \leq i < T} E\left[\log \left(\frac{W_{i+1}}{W_{i}} \right) \right] + \log(K+1) \quad (\text{Item 3}) \\ &\leq \frac{\eta}{1-\gamma} \cdot E[\tilde{\mathbb{U}}] + (e-2) \cdot \frac{\eta^{2}}{1-\gamma} \sum_{1 \leq i \leq T} \sum_{k} E\left[p_{ik} \cdot \widehat{U}_{ik}^{2} \right] \\ &\quad + \log(K+1) \quad (\text{Item 4 and 5}) \\ &\leq \frac{\eta}{1-\gamma} \cdot E[\tilde{\mathbb{U}}] + (e-2) \cdot \frac{\eta^{2}}{1-\gamma} T \cdot \frac{K+1}{K} \cdot \left(\frac{2K+1}{6} + \frac{\lambda^{2}}{\gamma} \right) \\ &\quad + \log(K+1). \quad (\text{Item 6}) \end{split}$$

Multiplying by $(1 - \gamma)$ and dividing by η , adding $\gamma \sup_x \mathbb{U}(x) + T\frac{\lambda}{K}$ to both sides and subtracting $E[\tilde{\mathbb{U}}]$, bounding $\sup_x \mathbb{U}(x) \leq T$, and $E[\tilde{\mathbb{U}}] \leq E[\mathbb{U}]$ (from Item 1), yields

$$\sup_{x} \mathbb{U}(x) - E[\mathbb{U}]$$

$$\leq \left(\gamma + \eta \cdot (e - 2)\frac{K + 1}{K} \cdot \left(\frac{2K + 1}{6} + \frac{\lambda^{2}}{\gamma}\right) + \frac{\lambda}{K}\right) \cdot T$$

$$+ \frac{\log(K + 1)}{\eta}.$$
(46)

This proves the first claim of the theorem. 8. *Optimizing tuning parameters* Suppose now that we choose the tuning parameters as follows:

$$\gamma = c_1 \cdot \left(\frac{\log(T)}{T}\right)^{1/3}, \qquad \eta = c_2 \cdot \gamma^2, \qquad K = c_3/\gamma.$$

Substituting we get

$$\begin{split} \sup_{x} \mathbb{U}(x) &- E[\mathbb{U}] \\ &\leq \left(\gamma + c_{2} \cdot \gamma^{2} \cdot (e - 2) \frac{K + 1}{K} \cdot \left(\frac{2c_{3}/\gamma + 1}{6} + \frac{\lambda^{2}}{\gamma}\right) + \lambda \cdot \gamma/c_{3}\right) \cdot T + \frac{\log(K + 1)}{c_{2} \cdot \gamma^{2}} \\ &= \log(T)^{1/3} T^{2/3} \cdot \left(c_{1} + (e - 2) \frac{K + 1}{K} \cdot c_{1}c_{2}\left(\frac{c_{3}}{3} + \lambda^{2} + \frac{\gamma}{6}\right) + \lambda \frac{c_{1}}{c_{3}} \\ &+ \frac{\log(T^{1/3}\log(T)^{-1/3}c_{3}/c_{1} + 1)}{c_{1}^{2}\log(T)}\right) \\ &= \log(T)^{1/3} T^{2/3} \cdot \left(c_{1} + (e - 2) \cdot c_{1}c_{2}\left(\frac{c_{3}}{3} + \lambda^{2}\right) + \lambda \frac{c_{1}}{c_{3}} + \frac{1}{3c_{1}^{2}} + o(1)\right). \end{split}$$

The second claim of the theorem follows.

REFERENCES

- ACHDDOU, JULIETTE, OLIVIER CAPPÉ, AND AURÉLIEN GARIVIER (2021): Fast Rate Learning in Stochastic First Price Bidding. Asian Conference on Machine Learning. [1077]
- AGRAWAL, RAJEEV (1995): "Sample Mean Based Index Policies by $\mathcal{O}(\log n)$ Regret for the Multi-Armed Bandit Problem," Advances in applied probability, 27 (4), 1054–1078. [1083]
- AUER, PETER, NICOLO CESA-BIANCHI, AND PAUL FISCHER (2002): "Finite-Time Analysis of the Multiarmed Bandit Problem," *Machine learning*, 47, 235–256. [1083]
- AUER, PETER, NICOLO CESA-BIANCHI, YOAV FREUND, AND ROBERT E. SCHAPIRE (2002): "The Nonstochastic Multiarmed Bandit Problem," *SIAM journal on computing*, 32 (1), 48–77. [1075,1082]
- BACHOC, FRANÇOIS, TOMMASO CESARI, ROBERTO COLOMBONI, AND ANDREA PAUDICE (2022a): "A Near-Optimal Algorithm for Univariate Zeroth-Order Budget Convex Optimization." [1087]
- (2022b): "Regret Analysis of Dyadic Search," arXiv preprint arXiv:2209.00885. [1087]
- BAILY, MARTIN N. (1978): "Some Aspects of Optimal Unemployment Insurance," *Journal of Public Economics*, 10 (3), 379–402. [1073,1076]
- CESA-BIANCHI, NICOLO, AND GÁBOR LUGOSI (2006): Prediction, Learning, and Games. Cambridge University Press. [1084,1099]
- CESA-BIANCHI, NICOLÒ, TOMMASO CESARI, ROBERTO COLOMBONI, FEDERICO FUSCO, AND STEFANO LEONARDI (2024a): "Bilateral Trade: A Regret Minimization Perspective," *Mathematics of Operations Research*, 49 (1), 171–203. [1077,1079,1080]
- CESA-BIANCHI, NICOLÒ, TOMMASO CESARI, ROBERTO COLOMBONI, FEDERICO FUSCO, AND STEFANO LEONARDI (2024b): "Regret Analysis of Bilateral Trade With a Smoothed Adversary," *Journal of Machine Learning Research*, 25 (234), 1–36. [1077,1080]
- CESA-BIANCHI, NICOLÒ, ROBERTO COLOMBONI, AND MAXIMILIAN KASY (2025): "Supplement to 'Adaptive Maximization of Social Welfare'," *Econometrica Supplemental Material*, 93, https://doi.org/10.3982/ ECTA22351. [1076]
- CHETTY, RAJ (2009): "Sufficient Statistics for Welfare Analysis: A Bridge Between Structural and Reduced-Form Methods," *Annual Review of Economics*, 1 (1), 451–488. [1073,1076,1094]
- COHEN, JESSICA, AND PASCALINE DUPAS (2010): "Free Distribution or Cost-Sharing? Evidence From a Randomized Malaria Prevention Experiment," *The Quarterly Journal of Economics*, 125 (1), 1–45. [1094]
- COVER, THOMAS (1965): "Behavior of Sequential Predictors of Binary Sequences," in *Proceedings of the 4th Prague Conference on Information Theory, Statistical Decision Functions, Random Processes.* Prague: Publishing House of the Czechoslovak Academy of Sciences, 263–272. [1077]

- DASKALAKIS, CONSTANTINOS, AND VASILIS SYRGKANIS (2022): "Learning in Auctions: Regret Is Hard, Envy Is Easy," *Games and Economic Behavior*. [1077]
- DEN BOER, ARNOUD V. (2015): "Dynamic Pricing and Learning: Historical Origins, Current Research, and New Directions," *Surveys in Operations Research and Management Science*. [1077,1094]
- DYNARSKI, SUSAN, LINDSAY PAGE, AND JUDITH SCOTT-CLAYTON (2023): "College Costs, Financial Aid, and Student Decisions," in *Handbook of the Economics of Education*, Vol. 7. Elsevier, 227–285. [1094]
- FENG, ZHE, GURU GURUGANESH, CHRISTOPHER LIAW, ARANYAK MEHTA, AND ABHISHEK SETHI (2021): "Convergence Analysis of No-Regret Bidding Algorithms in Repeated Auctions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 5399–5406. [1077]
- FENG, ZHE, CHARA PODIMATA, AND VASILIS SYRGKANIS (2018): "Learning to Bid Without Knowing Your Value," in *Proceedings of the 2018 ACM Conference on Economics and Computation*, 505–522. [1077]
- FUDENBERG, DREW, AND DAVID K. LEVINE (1995): "Consistency and Cautious Fictitious Play," Journal of Economic Dynamics and Control, 19 (5-7), 1065-1089. [1077]
- HAN, YANJUN, ZHENGYUAN ZHOU, AARON FLORES, ERIK ORDENTLICH, AND TSACHY WEISSMAN (2020): "Learning to Bid Optimally and Efficiently in Adversarial First-Price Auctions," arXiv preprint arXiv:2007. 04568. [1077]
- HAN, YANJUN, ZHENGYUAN ZHOU, AND TSACHY WEISSMAN (2020): "Optimal No-Regret Learning in Repeated First-Price Auctions," arXiv preprint arXiv:2003.09795. [1077]
- HANNAN, JAMES (1957): "Approximation to Bayes Risk in Repeated Play," Contributions to the Theory of Games, 3 (2), 97–139. [1077]
- HART, SERGIU, AND ANDREU MAS-COLELL (2000): "A Simple Adaptive Procedure Leading to Correlated Equilibrium," *Econometrica*, 68 (5), 1127–1150. [1077]
- (2001): "A General Class of Adaptive Strategies," Journal of Economic Theory, 98 (1), 26–54. [1077]
- KASY, MAXIMILIAN (2018): "Optimal Taxation and Insurance Using Machine Learning—Sufficient Statistics and Beyond," *Journal of Public Economics*, 167. [1095]
- KASY, MAXIMILIAN, AND ANJA SAUTMANN (2021): "Adaptive Treatment Assignment in Experiments for Policy Choice," *Econometrica*, 89 (1), 113–132. [1077]
- KLEINBERG, ROBERT D., AND FRANK THOMSON LEIGHTON (2003): "The Value of Knowing a Demand Curve: Bounds on Regret for Online Posted-Price Auctions," in *IEEE Symposium on Foundations of Computer Science*, 594–605. [1077,1079,1080]
- KOLUMBUS, YOAV, AND NOAM NISAN (2022): "Auctions Between Regret-Minimizing Agents," in Proceedings of the ACM Web Conference 2022, 100–111. [1077]
- LAI, TZE LEUNG (1987): "Adaptive Treatment Allocation and the Multi-Armed Bandit Problem," *The Annals of Statistics*, 1091–1114. [1083]
- LAI, TZE LEUNG, AND HERBERT ROBBINS (1985): "Asymptotically Efficient Adaptive Allocation Rules," Advances in applied mathematics, 6 (1), 4–22. [1073]
- LATTIMORE, TOR, AND CSABA SZEPESVÁRI (2020): *Bandit Algorithms*. Cambridge University Press. [1073, 1076,1085]
- LEE, KENNETH, EDWARD MIGUEL, AND CATHERINE WOLFRAM (2020): "Experimental Evidence on the Economics of Rural Electrification," *Journal of Political Economy*, 128 (4), 1523–1565. [1094]
- LITTLESTONE, NICK, AND MANFRED WARMUTH (1994): "The Weighted Majority Algorithm," Information and computation, 108 (2), 212–261. [1077]
- LUCIER, BRENDAN, SARATH PATTATHIL, ALEKSANDRS SLIVKINS, AND MENGXIAO ZHANG (2024): "Autobidders With Budget and ROI Constraints: Efficiency, Regret, and Pacing Dynamics," in *The Thirty Seventh Annual Conference on Learning Theory*, *PMLR*, 3642–3643. [1077]
- LUGOSI, GÁBOR, MIHALIS G. MARKAKIS, AND GERGELY NEU (2023): "On the Hardness of Learning From Censored and Nonstationary Demand," *INFORMS Journal on Optimization*. [1077]
- MILGROM, PAUL, AND ILYA SEGAL (2002): "Envelope Theorems for Arbitrary Choice Sets," *Econometrica*, 70 (2), 583–601. [1092,1094]
- MIRRLEES, JAMES A. (1971): "An Exploration in the Theory of Optimum Income Taxation," *The Review of Economic Studies*, 175–208. [1073,1076,1089]
- NEDELEC, THOMAS, CLÉMENT CALAUZÈNES, NOUREDDINE EL KAROUI, VIANNEY PERCHET et al. (2022): "Learning in Repeated Auctions," *Foundations and Trends*® *in Machine Learning*, 15 (3), 176–334. [1094]
- RAMSEY, FRANK P. (1927): "A Contribution to the Theory of Taxation," *The economic journal*, 37 (145), 47–61. [1073,1076,1089,1092]
- RUSSO, DANIEL (2020): "Simple Bayesian Algorithms for Best-Arm Identification," *Operations Research*, 68 (6), 1625–1647. [1077]
- RUSSO, DANIEL J., BENJAMIN VAN ROY, ABBAS KAZEROUNI, IAN OSBAND, AND ZHENG WEN (2018): "A Tutorial on Thompson Sampling," *Foundations and Trends in Machine Learning*, 11 (1), 1–96. [1095]

- SAEZ, EMMANUEL (2001): "Using Elasticities to Derive Optimal Income Tax Rates," *The Review of Economic Studies*, 68 (1), 205–229. [1073,1076,1089]
- (2002): "Optimal Income Transfer Programs: Intensive versus Extensive Labor Supply Responses," *The Quarterly Journal of Economics*, 117 (3), 1039–1073. [1089]
- SAEZ, EMMANUEL, AND STEFANIE STANTCHEVA (2016): "Generalized Social Welfare Weights for Optimal Tax Theory," *American Economic Review*, 106 (1), 24–45. [1089]
- SELDIN, YEVGENY, AND ALEKSANDRS SLIVKINS (2014): "One Practical Algorithm for Both Stochastic and Adversarial Bandits," in *International Conference on Machine Learning*, *PMLR*, 1287–1295. [1077]

SINANDER, LUDVIG (2022): "The Converse Envelope Theorem," *Econometrica*, 90 (6), 2795–2819. [1094]

SLIVKINS, ALEKSANDRS (2019): "Introduction to Multi-Armed Bandits," arXiv preprint arXiv:1904.07272. [1073,1076]

- THOMPSON, WILLIAM R. (1933): "On the Likelihood That One Unknown Probability Exceeds Another in View of the Evidence of Two Samples," *Biometrika*, 25 (3/4), 285–294. [1073,1095]
- VOVK, VOLODIMIR G. (1990): "Aggregating Strategies," in Proceedings of the 3rd AnnualWworkshop on Computational Learning Theory, 371–386. [1077]
- WEED, JONATHAN, VIANNEY PERCHET, AND PHILIPPE RIGOLLET (2016): "Online Learning in Repeated Auctions," in *Conference on Learning Theory, PMLR*, 1562–1583. [1077]
- WILLIAMS, CHRISTOPHER K. I., AND CARL E. RASMUSSEN (2006): Gaussian Processes for Machine Learning. MIT Press. [1095]
- ZIMMERT, JULIAN, AND YEVGENY SELDIN (2021): "Tsallis-INF: An Optimal Algorithm for Stochastic and Adversarial Bandits," *Journal of Machine Learning Research*, 22 (28), 1–49. [1077,1085]

Co-editor Guido W. Imbens handled this manuscript.

Manuscript received 15 October, 2023; final version accepted 22 March, 2025; available online 25 March, 2025.

The replication package for this paper is available at https://doi.org/10.5281/zenodo.15042114. The Journal checked the data and codes included in the package for their ability to reproduce the results in the paper and approved online appendices.