

# The Role of Transparency in Repeated First-Price Auctions with Unknown Valuations\*

Nicolò Cesa-Bianchi<sup>1,4</sup>, Tommaso Cesari<sup>2</sup>, Roberto Colomboni<sup>1,4</sup>,  
Federico Fusco<sup>3</sup>, and Stefano Leonardi<sup>3</sup>

<sup>1</sup>Università degli Studi di Milano, Milano, Italy

<sup>2</sup>University of Ottawa, Ottawa, Canada

<sup>3</sup>Sapienza Università di Roma, Roma, Italy

<sup>4</sup>Politecnico di Milano, Milano, Italy

## Abstract

We study the problem of regret minimization for a single bidder in a sequence of first-price auctions where the bidder discovers the item's value only if the auction is won. Our main contribution is a complete characterization, up to logarithmic factors, of the minimax regret in terms of the auction's *transparency*, which controls the amount of information on competing bids disclosed by the auctioneer at the end of each auction. Our results hold under different assumptions (stochastic, adversarial, and their smoothed variants) on the environment generating the bidder's valuations and competing bids. These minimax rates reveal how the interplay between transparency and the nature of the environment affects how fast one can learn to bid optimally in first-price auctions.

---

\*This is the full version of Cesa-Bianchi et al. [2024a]

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Overview of Our Results . . . . .	2
1.2	Technical Challenges . . . . .	3
1.3	Related Work . . . . .	4
<b>2</b>	<b>The Learning Model</b>	<b>6</b>
<b>3</b>	<b>The Stochastic i.i.d. Setting</b>	<b>7</b>
3.1	I.I.D. – Bandit Feedback . . . . .	7
3.2	I.I.D. – Semi-Transparent Feedback . . . . .	8
3.2.1	A $T^{2/3}$ upper bound for the i.i.d. environment . . . . .	8
3.2.2	A $T^{2/3}$ lower bound for the smooth i.i.d. environment . . . . .	10
3.3	I.I.D. – Transparent/Full Feedback . . . . .	12
3.3.1	A $\sqrt{T}$ upper bound for the i.i.d. environment . . . . .	12
3.3.2	A $\sqrt{T}$ lower bound for the i.i.d. environment . . . . .	15
<b>4</b>	<b>The Adversarial Setting</b>	<b>18</b>
4.1	Smooth – Bandit Feedback . . . . .	18
4.2	Smooth – Transparent Feedback . . . . .	19
4.3	The (Non-Smooth) Adversarial Model . . . . .	20
<b>5</b>	<b>Conclusion</b>	<b>21</b>
<b>A</b>	<b>Appendix</b>	<b>25</b>
A.1	Measure and Information-Theoretic Notation and Known Facts . . . . .	25
A.2	Missing Details of the proof of Theorem 3 . . . . .	26
A.3	Missing Details of the proof of Theorem 6 . . . . .	31

# 1 Introduction

The online advertising market has recently transitioned from second to first-price auctions. A remarkable example is Google AdSense’s move at the end of 2021 [Wong, 2021], following the switch made by Google AdManager and AdMob. Earlier examples include OpenX, AppNexus, Index Exchange, and Rubicon [Sluis, 2017]. To increase transparency in first-price auctions, some platforms (like AdManager) have a single bidding session for each available impression (unified bidding) and require all partners to share and receive bid data. After the first-price auction closes, bidders receive the minimum bid price that would have won them the impression [Bigler, 2019]. In practice, advertisers face two main sources of uncertainty in the bidding phase: they ignore the value of the competing bids and, crucially, ignore the actual value of the impression they are bidding on. Indeed, clicks and conversion rates—which are only measured *after* the auction is won and the ad is displayed—can vary wildly over time or highly correlate with competing bids. We remark that ignoring the value of the impression strongly affects the bidder’s utility: it may lead to overbidding for an impression of low value or, conversely, underbidding and losing a valuable one. To cope with this uncertainty, advertisers rely on auto-bidders that use the feedback provided in the auctions to learn good bidding strategies. We study the learning problem faced by a single bidder within the framework of regret minimization according to the following protocol:

---

## Online Bidding Protocol

---

**for**  $t = 1, 2, \dots, T$  **do**

Valuation  $V_t$  and competing bid  $M_t$  are privately generated

The learner posts a bid  $B_t$  and receives utility  $\text{Util}_t(B_t)$ :

$$\text{Util}_t(B_t) = (V_t - B_t)\mathbb{I}\{B_t \geq M_t\}$$

The learner observes some feedback  $Z_t$

---

The bidder has no initial information on the environment and seeks to learn the relevant features of the problem on the fly. The performance of a learning strategy for the bidder—also referred to as the learner—is measured in terms of the difference in total utility with respect to the best fixed bid. This difference is called *regret*, and the main goal is to design strategies with asymptotically vanishing time-averaged regret with respect to the best fixed-bid strategy or, equivalently, regret sublinear in the time horizon.

In this work, we are specifically interested in understanding how the “transparency” of the auctions—i.e., the amount of information on competing bids disclosed by the auctioneer *after* the auction takes place—affects the learning process. There is a clear tension regarding transparency: on the one hand, bidders want to receive as much information as possible about the environment to learn the competitor’s bidding strategies while revealing as little as possible about their (private) bids. On the other hand, the platform may not want to publicly reveal its revenue (i.e., the winning bid). Our investigation addresses both sides of the “transparency dilemma”. Our algorithmic results provide bidders with a toolbox of learning strategies to (optimally) exploit the various degrees of transparency, while the tightness of our results fully characterizes the impact of transparency on learnability. This complete picture allows platforms to make an informed decision in choosing their level of transparency, as it is in their interest to create a thriving environment for advertisers.

To model the level of transparency, we distinguish four natural types of feedback  $Z_t$ , specifying the conditions under which the highest competing bid  $M_t$  and the bidder’s valuation  $V_t$  are revealed to the bidder after each round  $t$ . In the transparent feedback setting,  $M_t$  is always observed after

	Stochastic i.i.d.		Adversarial	
	Smooth	General	Smooth	General
Full Feedback	Thm.5: $\Omega(\sqrt{T})$			Thm.8: $\Omega(T)$
Transparent		Thm.4: $O(\sqrt{T})$	Thm.7: $\tilde{O}(\sqrt{T})$	
Semi-Transparent	Thm.3: $\Omega(T^{2/3})$	Thm.2: $\tilde{O}(T^{2/3})$		
Bandit Feedback		Thm.1: $\Omega(T)$	Thm.6: $O(T^{2/3})$	

Table 1: Summary of our results. Rows correspond to feedback models while columns to environments. The minimax regret of every problem falls in one of the following three regimes:  $\tilde{\Theta}(\sqrt{T})$  (green),  $\tilde{\Theta}(T^{2/3})$  (yellow) and  $\tilde{\Theta}(T)$  (red).

the auction is concluded, while  $V_t$  is only known if the auction is won, i.e., when  $B_t \geq M_t$ . In the semi-transparent setting,  $M_t$  is only observed when the auction is lost. In other words, in the semi-transparent setting, the platform publicly reveals only the winning bid, whereas in the transparent setting, the platform reveals all bids. We also consider two extreme settings that provide two natural learning benchmarks: full feedback ( $M_t$  and  $V_t$  are always observed irrespective of the auction’s outcome) and bandit feedback ( $M_t$  is never observed while  $V_t$  is only observed by the winning bidder). Note that the learner can compute the value of the utility  $\text{Util}_t(B_t)$  at time  $t$  with any type of feedback, including bandit feedback. In this paper, we characterize the learner’s minimax regret not only with respect to the degree of transparency of the auction but also with respect to the nature of the process generating the sequence of pairs  $(V_t, M_t)$ . In particular, we consider four types of environments: stochastic i.i.d., adversarial, and their smooth versions (see Section 1.3 for a discussion about smoothness, and Section 2 for the formal definition).

## 1.1 Overview of Our Results

We report here an overview of our results (see also Table 1). For simplicity, we often hide the logarithmic factors with the  $\tilde{O}$  notation.

### Stochastic i.i.d. settings

- In both the full and transparent feedback models, the minimax regret is of order  $\sqrt{T}$  (Theorems 4 and 5), and adding the smoothness requirement leaves this rate unchanged.
- In the semi-transparent feedback model, the minimax regret is of order  $T^{2/3}$  (Theorems 2 and 3). Also in this case, adding the smoothness requirement leaves this rate unchanged.
- In the bandit feedback model, smoothness is crucial for sublinear regret (Theorem 1). In particular, smoothness implies a minimax regret of  $T^{2/3}$  (this is obtained by combining the upper bound in Theorem 6 and the lower bound in Theorem 3).

### Adversarial settings

- Without smoothness, sublinear regret cannot be achieved, even with full feedback (Theorem 8).
- In both the full and transparent feedback model, the minimax regret in a smooth environment is of order  $\sqrt{T}$  (combining the lower bound in Theorem 5 and the upper bound in Theorem 7).
- Both with semi-transparent and bandit feedback, the minimax regret in a smooth environment is of order  $T^{2/3}$  (combining the lower bound in Theorem 3 and the upper bound in Theorem 6).

Interestingly, the minimax regret rates for first-price auctions mirror the allowed regret regimes in finite partial monitoring games [Bartók et al., 2014] and online learning with feedback graphs [Alon

et al., 2017]. This is somehow surprising, as it has been shown in Lattimore [2022] that games with continuous outcome/action spaces allow for a much larger set of regret rates—see also Cesa-Bianchi et al. [2023, 2024b], Bolić et al. [2024], Bernasconi et al. [2024].

Table 1 reveals some interesting properties of the learnability of the problem: full feedback and transparent feedback are essentially equivalent, while semi-transparent feedback and bandit feedback differ only in the stochastic i.i.d. setting. Qualitatively, this tells the platform that disclosing all bids (instead of only the winning one) drastically improves the learnability of the problem (green vs. yellow entries in Table 1). Besides, revealing at least the winning bid avoids some pathological behavior (yellow entries vs. red entry for the general i.i.d. environment with bandit feedback). Moreover, while smoothness is key for learning in the adversarial setting, in the stochastic case smoothness is only relevant for bandit feedback.

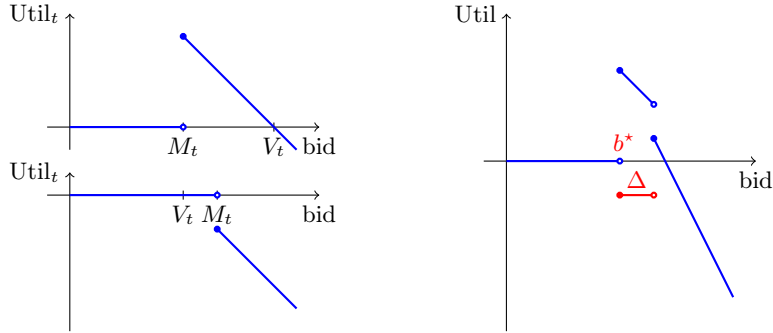


Figure 1: The utility function is generally neither Lipschitz nor continuous. If  $M_t \leq V_t$  (top left plot), then  $\text{Util}_t$  is upper-semi continuous and one-sided Lipschitz; conversely, if  $M_t \geq V_t$  (bottom left plot), then  $\text{Util}_t$  is still one-sided Lipschitz—from the other side—and lower-semi continuous. Summing up the two types of utilities results in a total utility that may be neither one-sided Lipschitz nor semi-continuous (right plot, where the two utility functions of the other two plots are summed up. There,  $b^*$  is the optimal bid and  $\Delta$  is the neighborhood of  $b^*$  where the total utility is “good enough”).

## 1.2 Technical Challenges

**The utility function.** The utilities  $\text{Util}_t(b) = (V_t - b)\mathbb{I}\{M_t \leq b\}$  are defined over a continuous decision space  $[0, 1]$  and are neither Lipschitz nor continuous, see Figure 1. Actually, even weaker properties, i.e., that the expected cumulative reward  $b \mapsto \sum_{t \in [T]} \mathbb{E}[\text{Util}_t(b)]$  is one-sided Lipschitz or semi-continuous, do not hold in general. We address this problem by developing techniques designed to control the approximation error incurred when discretizing the bidding space. This is a non-trivial problem without regularity assumption, as the neighborhood of the optimal bid where the total utility is “good enough” can be arbitrarily small in general (see the red interval  $\Delta$  in the rightmost plot of Figure 1). In the stochastic i.i.d. setting, the approximation error is controlled by building a sample-based *non-uniform* grid of candidate bids, which can be of independent interest. This allows us to estimate the distribution of the competing bids uniformly over the subintervals of  $[0, 1]$ . In the adversarial setting, instead, we use the smoothness assumption to guarantee that the expected utility is Lipschitz. In this case, the approximation error is controlled using a uniform grid with an appropriate grid-size (Lemma 4).

**The feedback models.** Our feedback models interpolate between bandit (only the bidder’s utility is observed) and full feedback ( $V_t$  and  $M_t$  are always observed). In the stochastic i.i.d. case, the different levels of transparency are crucial to the process of building the non-uniform grids used to

control the discretization error. In the adversarial case, when there are only  $K$  allowed bids, the optimal rates are of order  $\sqrt{T \ln K}$  and  $\sqrt{KT}$  under full and bandit feedback, respectively. While the semi-transparent feedback is not enough to improve on the bandit rate, the transparent one can be exploited via a more sophisticated approach. To this end, we design an algorithm, Exp3.FPA, enjoying the full feedback regret rate of order  $\sqrt{T \ln K}$  while only relying on the weaker transparent feedback.

**Lower bounds.** The linear lower bounds (Theorems 1 and 8) exploit a “needle in a haystack” phenomenon, where there is a hidden optimal bid  $b^*$  in the  $[0, 1]$  interval and the learner has no way of finding  $b^*$  using the feedback it has access to. This is indeed the case in the non-smooth adversarial full-feedback setting and in the non-smooth i.i.d. bandit setting. To prove the remaining lower bounds, we design careful embeddings of known hard instances into our framework. In particular, in Theorem 5 we embed the hard instance for prediction with two experts and in Theorem 3 the hard instance for  $K = \Theta(T^{1/3})$  bandits.

### 1.3 Related Work

**Transparency in first-price auctions.** The role of transparency in repeated first-price auctions has been investigated by Bergemann and Hörner [2018], but mostly from a game-theoretic viewpoint. In particular, they study the impact of the feedback policy on the bidders’ strategy and show how disclosing the bids at the end of each round affects the equilibria of a bidding game with infinite horizon. In contrast, we want to characterize the impact of different amounts of feedback (or degrees of transparency) on the learner’s regret, which is measured against the optimal fixed bid in hindsight.

**Auctions with unknown valuations.** Although the problem of regret minimization in first-price auctions has been studied before, only a few papers consider the natural setting of unknown valuations. Feng et al. [2018] introduce a general framework for the study of regret in auctions where a bidder’s valuation is only observed when the auction is won. In the special case of first-price auctions, their setting is equivalent to our transparent feedback when the sequence of pairs  $(V_t, M_t)$  is adversarially generated. Following a parameterization introduced by Weed et al. [2016], Feng et al. [2018] provide a  $O(\sqrt{T \ln \max\{\Delta_0^{-1}, T\}})$  regret bound, where  $\Delta_0 = \min_{t < t'} |M_t - M_{t'}|$  is controlled by the environment. In the stochastic i.i.d. case, their results translate into *distribution-dependent* guarantees that do not translate into a worst-case sublinear bound (we obtain a  $\sqrt{T}$  rate). In the adversarial case, their guarantees are still linear in the worst-case (we obtain  $\sqrt{T}$  bounds by leveraging the smoothness assumption). Achddou et al. [2021] consider a stochastic i.i.d. setting with the additional assumption that  $V_t$  and  $M_t$  are independent. Their main result is a bidding algorithm with *distribution-dependent* regret rates (of order  $T^{1/3+\varepsilon}$  or  $\sqrt{T}$ , depending on the assumptions on the underlying distribution) in the transparent setting. Again, this result is not comparable to ours because of the independence assumption and the distribution-dependent rates (which do not allow to recover our minimax rates). Other works consider regret minimization in repeated second-price auctions with unknown valuations. Dikkala and Tardos [2013] investigate a repeated bidding setting, but do not consider regret minimization. Weed et al. [2016] derive regret bounds for the case when  $M_t$  are adversarially generated, while  $V_t$  are stochastically or adversarially generated and the feedback is transparent.

**First-price auctions with known valuations.** Considerably more works study first price auctions when the valuation  $V_t$  is known to the bidder at the beginning of each round  $t$ . Note that

these results are not directly comparable to ours. Balseiro et al. [2019] look at the case when the  $V_t$  are adversarial and the  $M_t$  are either stochastic i.i.d. or adversarial. In the bandit feedback case (when  $M_t$  is never observed), they show that the minimax regret is  $\tilde{\Theta}(T^{2/3})$  in the stochastic case and  $\tilde{\Theta}(T^{3/4})$  in the adversarial case. Han et al. [2020b] prove a  $\tilde{O}(\sqrt{T})$  regret bound in the semi-transparent setting ( $M_t$  observed only when the auction is lost) with adversarial valuations and stochastic bids. Han et al. [2020a] focus on the adversarial case, when  $V_t$  and  $M_t$  are both generated adversarially. They prove a  $\tilde{O}(\sqrt{T})$  regret bound in the full feedback setting ( $M_t$  always observed) when the regret is defined with respect to all Lipschitz shading policies. This setup is extended in Zhang et al. [2022] where the authors consider the case in which the bidder is provided access to hints before each auction. Zhang et al. [2021] also studied the full information feedback setting and design a space-efficient variant of the algorithm proposed by Han et al. [2020a]. Badanidiyuru et al. [2023] introduce a contextual model in which  $V_t$  is adversarial and  $M_t = \langle \theta, x_t \rangle + \varepsilon_t$  where  $x_t \in \mathbb{R}^d$  is contextual information available at the beginning of each round  $t$ ,  $\theta \in \mathbb{R}^d$  is an unknown parameter, and  $\varepsilon_t$  is drawn from an unknown log-concave distribution. They study regret in bandit and full feedback settings.

**Dynamics in first-price auctions.** A different thread of research is concerned with the convergence property of the regret minimization dynamics in first-price auctions (or, more specifically, with the learning dynamics of mean-based regret minimization algorithms). Feldman et al. [2016] show that with continuous bid levels, coarse-correlated equilibria exist whose revenue is below the second price. Feng et al. [2021] prove that regret minimizing bidders converge to a Bayesian Nash equilibrium in a first-price auctions when bidder values are drawn i.i.d. from a uniform distribution on  $[0, 1]$ . Kolumbus and Nisan [2022] show that if two bidders with finitely many bid values converge, then the equilibrium revenue of the bidder with the highest valuation is the second price. Deng et al. [2022] characterize the equilibria of the learning dynamics depending on the number of bidders with the highest valuation. Their characterization is for both time-average and last-iterate convergence.

**Smoothed adversary.** Smoothed analysis of algorithms, originally introduced by Spielman and Teng [2004] and later formalized for online learning by Rakhlin et al. [2011], Haghtalab et al. [2020], is a known approach to the analysis of algorithms in which the instances at every round are generated from a distribution that is not too concentrated. Recent works on the smoothed analysis of online learning algorithms include Kannan et al. [2018], Haghtalab et al. [2020, 2022], Block et al. [2022], Durvasula et al. [2023], Cesa-Bianchi et al. [2023, 2024b, 2021, 2024c], Bolić et al. [2024].

**Online learning in metric spaces.** Our problem is related to online learning in metric spaces [Kleinberg et al., 2019], where the action space is endowed with a metric and the losses are induced by a sequence of Lipschitz functions defined onto it. Tight regret bounds are known, parameterized by some notion of dimension of the metric space, in both the full and the bandit models. The simple structure of our action space ( $[0, 1]$  with the Euclidean distance) allows us to obtain tight bounds by either using a uniform grid (Theorems 6 and 7) or sample-based grids (Theorems 2 and 4), without resorting to the more elaborate techniques that characterize this line of research, e.g., zooming (which is typically used in the bandit feedback model to account for the lack of feedback). Also related to our model is the study of piecewise and regular Lipschitz functions [Balcan et al., 2018, Sharma et al., 2020, Duetting et al., 2023]. In particular, Lemma 1 and Theorem 3 in Balcan et al. [2018] imply our Theorem 6 in the special case of independent processes.\*

---

\*Combining the second part of their Lemma 1 with their Theorem 3 to lift independence gives void guarantees in the general case (note that there is a typo in the statement of their Lemma 1: as it can be seen in the proof, the

## 2 The Learning Model

We introduce formally the repeated bidding problem in first-price auctions. At each time step  $t$ , a new item arrives for sale, for which the learner holds some unknown valuation  $V_t \in [0, 1]$ . The learner bids some  $B_t \in [0, 1]$  and, at the same time, a set of competitors bid for the same object. We denote their highest competing bid by  $M_t \in [0, 1]$ . The learner gets the item at cost  $B_t$  if it wins the auction (i.e., if  $B_t \geq M_t$ ), and does not get it otherwise. Then, the learner observes some feedback  $Z_t$  and gains utility  $\text{Util}_t(B_t)$ , where, for all  $b \in [0, 1]$ ,  $\text{Util}_t(b) = (V_t - b)\mathbb{I}\{b \geq M_t\}$  (see the Protocol in Section 1). Crucially, at time  $t$  the learner does not know its valuation  $V_t$  for the item before bidding, implying that its bid  $B_t$  only depends on its past observations  $Z_1, \dots, Z_{t-1}$  (and, possibly, some internal randomization). The goal of the learner is to design a learning algorithm  $\mathcal{A}$  that maximizes its utility. More precisely, we measure the performance of an algorithm  $\mathcal{A}$  by its *regret*  $R_T(\mathcal{A})$  against the worst environment  $\mathcal{S}$  in a certain class  $\Xi$ :  $R_T(\mathcal{A}) = \sup_{\mathcal{S} \in \Xi} R_T(\mathcal{A}, \mathcal{S})$ , where

$$R_T(\mathcal{A}, \mathcal{S}) = \sup_{b \in [0, 1]} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(b) - \sum_{t=1}^T \text{Util}_t(B_t) \right].$$

The expectation in the previous display is taken with respect to the randomness of the algorithm  $\mathcal{A}$  which selects  $B_t$ , and (possibly) the randomness of the environment  $\mathcal{S}$  generating the  $(V_t, M_t)$  pairs.

**The environments.** In this paper we consider both stochastic i.i.d. and adversarial environments.

- Stochastic i.i.d.: The pairs  $(V_1, M_1), (V_2, M_2), \dots$  are a stochastic i.i.d. process.
- Adversarial: The sequence  $(V_1, M_1), (V_2, M_2), \dots$  is generated by an oblivious adversary.

Following previous works in online learning (see Section 1.3), we also study versions of the above environments that are constrained to generate the sequence of  $(V_t, M_t)$  values using distributions that are “not too concentrated”. To this end, we introduce the notion of smooth distributions.

**Definition 1** (Haghtalab et al. [2021]). *Let  $\mathcal{X}$  be a domain that supports a uniform distribution  $\nu$ . A measure  $\mu$  on  $\mathcal{X}$  is said to be  $\sigma$ -smooth if for all measurable subsets  $A \subseteq \mathcal{X}$ , we have  $\mu(A) \leq \frac{\nu(A)}{\sigma}$ .*

We thus also consider the following two types of environments.

- The  $\sigma$ -smooth stochastic i.i.d. environment, which is a stochastic i.i.d. environment where the common distribution of all pairs  $(V_1, M_1), (V_2, M_2), \dots$  is  $\sigma$ -smooth.
- The  $\sigma$ -smooth adversarial setting, where the pairs  $(V_1, M_1), \dots$  form a stochastic process such that, for each  $t$ , the distribution of the pair  $(V_t, M_t)$  is  $\sigma$ -smooth.

**The feedback.** After describing the environments that we study, we now specify the types of feedback the learner receives at the end of each round, from the richest to the least informative.

- Full feedback. The learner observes its valuation and the highest competing bid:  $Z_t = (V_t, M_t)$ .
- Transparent feedback. The learner always observes  $M_t$ , but  $V_t$  is only revealed if it gets the item:  $Z_t$  is equal to  $(\star, M_t)$  if  $B_t < M_t$  and to  $(V_t, M_t)$  otherwise.
- Semi-transparent feedback<sup>†</sup>. The learner observes  $V_t$  if it gets the item and  $M_t$  otherwise:  $Z_t$  is equal to  $(\star, M_t)$  if  $B_t < M_t$  and to  $(V_t, \star)$  otherwise.
- The bandit feedback<sup>‡</sup>. The learner observes  $V_t$  if it gets the item and the symbol  $\star$  otherwise:  $Z_t$  is  $\star$  if  $B_t < M_t$  and to  $V_t$  otherwise.

correct result is  $k = P \cdot \mathcal{O}(M \cdot \kappa \cdot w + \sqrt{M \log(P/\zeta)})$  and, without assuming independence,  $P = T$  in our setting).

<sup>†</sup>This feedback is similar to the winner-only feedback in Han et al. [2020b].

<sup>‡</sup>We call this the bandit feedback because it is equivalent to receiving  $\text{Util}_t(B_t)$  (with the extra information  $\star$  to distinguish between losing the item and winning it with  $V_t = B_t$ , which does not affect regret guarantees).



### 3 The Stochastic i.i.d. Setting

In this section, we investigate the problem of repeated bidding in first-price auctions with unknown valuations, when the pairs of valuations and highest competing bids are drawn i.i.d. from a fixed but unknown distribution. We start by proving in Section 3.1 that it is impossible to achieve sublinear regret under the bandit feedback model without any assumption on the distribution of the environment. Then, in Section 3.2, we give matching upper and lower bounds of order  $T^{2/3}$  in the semi-transparent feedback model. Notably, the lower bound holds for smooth distributions, while the upper bound works for any (possibly non-smooth) distributions. Finally, in Section 3.3 we prove that both the full and transparent feedback yield the same minimax regret regime of order  $\sqrt{T}$ , regardless of the regularity of the distribution.

#### 3.1 I.I.D. – Bandit Feedback

In the bandit feedback model, at each time step, the learner observes the valuation  $V_t$  (and nothing else) when the auction is won, and observes nothing when the auction is lost. The crucial difference with the other (richer) types of feedback is the amount of information received about  $M_t$ , which, in the bandit case, is just the relative position with respect to  $B_t$  (i.e., whether  $M_t \leq B_t$  or  $B_t < M_t$ ). This allows to hide in the interval  $[0, 1]$  an optimal bid  $b^*$  which the learner cannot uncover over a finite time horizon. Following this idea, a difficult environment should randomize between two scenarios: a good scenario with large value  $V_t = 1$  and  $M_t$  slightly smaller than  $b^*$  and a bad one with poor value  $V_t = 0$  and  $M_t$  slightly larger than  $b^*$ . Then, to avoid suffering linear regret, the learner has to find this tiny interval around  $b^*$  (the “needle in a haystack”).

**Theorem 1.** *Consider the problem of repeated bidding in first-price auctions in a stochastic i.i.d. environment with bandit feedback. Then, any learning algorithm  $\mathcal{A}$  satisfies  $R_T(\mathcal{A}) \geq \frac{1}{13}T$ .*

*Proof.* We construct a randomized i.i.d. environment  $\mathcal{S}$ , such that any deterministic algorithm  $\mathcal{A}$  suffers linear regret against it, and then apply Yao’s minimax principle to conclude the proof. The randomized environment is simple: before starting the sequence, a uniform seed  $b^*$  is drawn uniformly at random in  $(1/3, 1/2 - \varepsilon)$ , where  $\varepsilon$  is a small parameter we set later. Then, the i.i.d. sequence  $(V_1, M_1), (V_2, M_2), \dots$  is drawn as follows: at each time step  $t$  with probability  $1/2$  we have  $(V_t, M_t) = (1, b^*)$ , otherwise  $(0, b^* + \varepsilon)$ . The best bid in hindsight,  $b^*$ , yields an overall expected utility of  $\frac{T}{2}(1 - b^*)$ , which is at least  $T/4$ , as  $b^*$  belongs to  $(1/3, 1/2)$ .

We now upper bound the utility achievable by any deterministic algorithm  $\mathcal{A}$  against  $\mathcal{S}$ . Fix any such algorithm, and consider its bids against any environment that selects the valuations  $V_t$  to be either 0 or 1 (as the one we just constructed). At each time step, the feedback that  $\mathcal{A}$  receives is 0, 1 or  $\star$  (when the item is allocated to one of the competitors), so that the history of the bids posted by  $\mathcal{A}$  is naturally described by a ternary decision tree of height  $T$ , where each level corresponds to a time step and any node to a bid. Crucially, the leaves of this tree are finite (at most  $3^T$ ), which means that the algorithm  $\mathcal{A}$  only posts bids in a finite subset  $N$  of  $[0, 1]$ . Now, let  $\varepsilon = 3^{-2T}/12$ ; we have that, with probability at least  $1 - 6N\varepsilon/(1-6\varepsilon) \geq 1 - e^{-T}$ , the set  $[b^*, b^* + \varepsilon]$  does not intersect  $N$ . Note: the randomness is with respect to the uniform seed  $b^*$  drawn by  $\mathcal{S}$ , while the bound on the probability holds independently to the choice of the deterministic algorithm  $\mathcal{A}$ .

The total utility of  $\mathcal{A}$  when  $[b^*, b^* + \varepsilon]$  does not intersect  $N$  is easy to analyze: every time that  $\mathcal{A}$  posts bids smaller than  $b^*$ , then it never wins the item (zero utility). Instead, if it posts bids larger than  $b^* + \varepsilon$ , then it always gets the item (whose average value is  $1/2$ ), paying at least  $b^* + \varepsilon \geq 1/3$ . Putting these two cases together, we have proved that at each time step the expected utility earned by the learner is at most  $1/6 = 1/2 - 1/3$ , when  $[b^*, b^* + \varepsilon] \cap N = \emptyset$  (which happens with probability

---

**COLLECT BIDS**


---

```

1: input: Time horizon  $T_0$ 
2:  $X_0 \leftarrow 0$  and  $M^{(0)} \leftarrow 0$ 
3: for each round  $t = 1, 2, \dots, T_0$  do
4:   Post bid  $B_t = 0$  and observe the highest competing bid  $M_t$ 
5:   Sort the observed highest competing bids in increasing order:  $M^{(1)} \leq M^{(2)} \leq \dots \leq M^{(T_0)}$ 
6:   if  $M^{(T_0)} = 0$  then return candidate bid  $X_0$ 
7:   for  $i = 1, 2, \dots$  do
8:      $j_{i-1}^* \leftarrow \max\{j \in \{0, \dots, T_0\} \mid X_{i-1} = M^{(j)}\}$ 
9:      $j_i \leftarrow \min\{j_{i-1}^* + \lceil \sqrt{T_0} \rceil, T_0\}$ ,  $X_i \leftarrow M^{(j_i)}$ 
10:    if  $j_i = T_0$  then let  $K \leftarrow i$  and break;
11: return Candidate bids  $X_0, X_1, X_2, \dots, X_K$ 

```

---

at least  $1 - e^{-T}$ ). Finally, by combining the lower bound on the performance of  $b^*$  with the upper bound on the expected utility of the learner, we get  $R_T(\mathcal{A}, \mathcal{S}) \geq (1 - e^{-T})(T/4 - T/6) \geq T/13$ .  $\square$

### 3.2 I.I.D. – Semi-Transparent Feedback

In this section, we prove two results settling the minimax regret for the semi-transparent feedback where the environment is i.i.d. (and, possibly, smooth). First, we construct a learning algorithm, COLLECTING BANDIT, achieving  $T^{2/3}$  regret against any i.i.d. environment. Then, we complement it with a lower bound of the same order (up to log terms) obtained even in a smooth i.i.d. environment.

#### 3.2.1 A $T^{2/3}$ Upper Bound for the i.i.d. Environment

Our learning algorithm COLLECTING BANDIT is composed of two phases. First, for  $T_0 = \Theta(T^{2/3})$  rounds, it collects samples from the highest competing bid random variables  $M_1, M_2, \dots, M_{T_0}$  by posting dummy bids  $B_1 = B_2 = \dots = B_{T_0} = 0$ . Among these values (plus the value  $X_0 = 0$ ), the algorithm selects  $\Theta(\sqrt{T_0})$  candidate bids according to their ordering, in such a way that the empirical frequencies of bids  $M_1, M_2, \dots, M_{T_0}$  landing strictly in between two consecutive selected values are at most  $\Theta(1/\sqrt{T_0})$  (see the pseudocode of COLLECT BIDS for details). Second, for the remaining time steps, it runs any bandit algorithm, using as candidate bids the ones collected in the first phase (see COLLECTING BANDIT for details). Note that, in this second phase, the (less informative) bandit feedback would be enough to run the algorithm: the additional information provided by the semi-transparent feedback is only exploited in the initial “collecting bids” phase. As a first step, we state a simple concentration result pertaining the i.i.d. process  $M, M_1, M_2, \dots, M_{T_0}$ , for  $T_0 \in \mathbb{N}$ . If  $\mathcal{I}$  is the family of all the subintervals of  $[0, 1]$  and  $\delta \in (0, 1)$ , we define

$$\mathcal{E}_\delta^{T_0} = \bigcap_{I \in \mathcal{I}} \left\{ \left| \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{I}\{M_t \in I\} - \mathbb{P}[M \in I] \right| < 8\sqrt{\frac{\ln(1/\delta)}{T_0}} \right\}.$$

The family  $\mathcal{I}$  of all the subintervals of  $[0, 1]$  has VC dimension 2 (see, e.g., Mitzenmacher and Upfal [2017, Chapter 14.2]). Therefore,  $\mathcal{E}_\delta^{T_0}$  is realized with probability at least  $1 - \delta$ , via standard sample complexity bound for  $\varepsilon$ -samples (see, e.g., Mitzenmacher and Upfal [2017, Theorem 14.15]). This is summarized in the following lemma.

**Lemma 1.** *For every  $T_0 \in \mathbb{N}$  and  $\delta \in (0, 1)$ , we have  $\mathbb{P}[\mathcal{E}_\delta^{T_0}] \geq 1 - \delta$ .*

For the sake of readability, we introduce the following notation:

**Notation 1.** Let  $\mathcal{X} = \{x_0, \dots, x_K\}$  be any grid with  $0 = x_0 < x_1 < \dots < x_K \leq 1$ , we denote by  $k_{\mathcal{X}}: [0, 1] \rightarrow \{0, 1, \dots, K\}$  the function that maps each  $b \in [0, 1]$  to the unique  $k$  such that  $b \in [x_k, x_{k+1})$ , with the convention that  $x_{K+1} = 2$ .

We now prove a lemma that allows us to control the expected cumulative utility of any bid in  $[0, 1]$  with that of the best bid in a discretization (without relying on any smoothness assumption).

**Lemma 2.** Consider any finite grid  $\mathcal{X} = \{x_0, \dots, x_K\}$ , with  $0 = x_0 < x_1 < \dots < x_K \leq 1$ , and assume that the process  $M, M_1, M_2, \dots$  of the highest competing bids form an i.i.d. sequence. For all  $b \in [0, 1]$  and  $T_0, T_1 \in \mathbb{N}$  with  $T_0 < T_1$ ,  $\mathbb{E} \left[ \sum_{t=T_0+1}^{T_1} \text{Util}_t(b) \right]$  is at most

$$\mathbb{E} \left[ \sum_{t=T_0+1}^{T_1} \text{Util}_t(x_{k_{\mathcal{X}}(b)}) \right] + (T_1 - T_0) \mathbb{P}[x_{k_{\mathcal{X}}(b)} < M < x_{k_{\mathcal{X}}(b)+1}].$$

*Proof.* Fix any  $b \in [0, 1]$ ,  $T_0, T_1 \in \mathbb{N}$  with  $T_0 < T_1$ , and a time step  $t \in \{T_0 + 1, \dots, T_1\}$ . Then

$$\begin{aligned} \mathbb{E}[\text{Util}_t(b)] &= \mathbb{E}[(V_t - b) \mathbb{I}\{b \geq M_t\}] \\ &\leq \mathbb{E}[(V_t - x_{k_{\mathcal{X}}(b)}) (\mathbb{I}\{x_{k_{\mathcal{X}}(b)} \geq M_t\} + \mathbb{I}\{b \geq M_t > x_{k_{\mathcal{X}}(b)}\})] \\ &\leq \mathbb{E}[\text{Util}_t(x_{k_{\mathcal{X}}(b)})] + \mathbb{P}[x_{k_{\mathcal{X}}(b)} < M_t \leq b] \\ &\leq \mathbb{E}[\text{Util}_t(x_{k_{\mathcal{X}}(b)})] + \mathbb{P}[x_{k_{\mathcal{X}}(b)} < M_t < x_{k_{\mathcal{X}}(b)+1}]. \end{aligned}$$

Summing over all times  $t$  and recalling that  $M_t$  and  $M$  share the same distribution, yields the conclusion.  $\square$

As a corollary of Lemmas 1 and 2 we obtain a similar discretization error guarantee when the grid of points  $\mathcal{X}$  is random.

**Lemma 3.** Fix any  $T_0 \in \mathbb{N}$  and  $\delta \in (0, 1)$ . Let  $\mathcal{X} = \{X_0, \dots, X_K\}$  be a random set containing a random number  $K$  of points satisfying  $0 = X_0 < X_1 < \dots < X_K \leq 1$ . Assume that the random variables  $K, X_0, X_1, \dots, X_{K+1}$  are  $\mathcal{H}_{T_0}$ -measurable, where  $\mathcal{H}_{T_0}$  is the history up to and including time  $T_0$ . Assume that the process  $(V_1, M_1), (V_2, M_2), \dots$  of the valuations/highest competing bids form an i.i.d. sequence. Then, for all  $b \in [0, 1]$  and  $T_1 \in \mathbb{N}$  with  $T_1 > T_0$ , we have:

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=T_0+1}^{T_1} \text{Util}_t(b) \right] &\leq \mathbb{E} \left[ \sum_{t=T_0+1}^{T_1} \text{Util}_t(X_{k_{\mathcal{X}}(b)}) \right] \\ &+ (T_1 - T_0) \left( \frac{1}{T_0} \sum_{t=1}^{T_0} \mathbb{P}[X_{k_{\mathcal{X}}(b)} < M_t < X_{k_{\mathcal{X}}(b)+1}] + 8\sqrt{\frac{\ln(1/\delta)}{T_0}} + \delta \right). \end{aligned}$$

We are now ready to present the main theorem of this section.

**Theorem 2.** Consider the problem of repeated bidding in first-price auctions in a stochastic i.i.d. environment with semi-transparent feedback. Then there exists a learning algorithm  $\mathcal{A}$  such that

$$R_T(\mathcal{A}) \leq 16(13 + \sqrt{\ln T})T^{2/3}.$$

*Proof.* We prove that COLLECTING BANDIT yields the desired bound when its learning routine  $\tilde{\mathcal{A}}$  is (a rescaled version of) MOSS [Audibert and Bubeck, 2009]: since MOSS is designed to run with gains in  $[0, 1]$  while the utilities we observe are in  $[-1, 1]$ , we first apply the reward transformation

- 1: **input:** Time horizon  $T$ , bandit algorithm  $\tilde{\mathcal{A}}$  for gains in  $[-1, 1]$
  - 2:  $T_0 \leftarrow \lceil T^{2/3} \rceil$
  - 3: Run COLLECT BIDS with horizon  $T_0$  and obtain  $X_0, X_1, \dots, X_K$
  - 4: Initialize  $\tilde{\mathcal{A}}$  on  $K + 1$  actions (one for each candidate bid  $X_i$ ) and  $T - T_0$  as time horizon
  - 5: **for** each round  $t = T_0 + 1, T_0 + 2, \dots, T$  **do**
  - 6:     Receive from  $\tilde{\mathcal{A}}$  the bid  $B_t = X_{I_t}$  for some  $I_t \in \{0, 1, \dots, K\}$
  - 7:     Post bid  $B_t$  and observe feedback  $Z_t$
  - 8:     Reconstruct  $\text{Util}_t(B_t)$  from  $Z_t$  and feed it to  $\tilde{\mathcal{A}}$
- 

$x \mapsto \frac{x+1}{2}$  to the observed utilities. This costs a multiplicative factor of 2 on the regret guarantees of MOSS. Leveraging the fact that the empirical frequency between two consecutive  $X_k$  and  $X_{k+1}$  generated by COLLECT BIDS is at most  $2/\sqrt{T_0}$  by design and applying Lemma 3 with  $T_1 = T$  to the random variables  $X_0, X_1, \dots, X_K$ , we get, for all  $b \in [0, 1]$ , that

$$\mathbb{E} \left[ \sum_{t=T_0+1}^T \text{Util}_t(b) \right] \leq \mathbb{E} \left[ \sum_{t=T_0+1}^T \text{Util}_t(X_{k_{\mathcal{X}}(b)}) \right] + (T - T_0) \left( \frac{2}{\sqrt{T_0}} + 8\sqrt{\frac{\ln(1/\delta)}{T_0}} + \delta \right) = (\star).$$

Now, applying the tower rule to the expectation on the right-hand side conditioning to the history  $\mathcal{H}_{T_0}$  up to time  $T_0$ , we can use the fact that the regret of the rescaled version of MOSS is upper bounded by  $98\sqrt{(K+1)(T-T_0)}$  and the number of points  $K+1$  collected by COLLECT BIDS is at most  $\sqrt{T_0} + 1$  to obtain

$$(\star) \leq \mathbb{E} \left[ \sum_{t=T_0+1}^T \text{Util}_t(B_t) \right] + 98\sqrt{(\sqrt{T_0} + 1)(T - T_0)} + (T - T_0) \left( \frac{2}{\sqrt{T_0}} + 8\sqrt{\frac{\ln(1/\delta)}{T_0}} + \delta \right).$$

Finally, tuning  $\delta = 1/T_0$ , upper bounding the cumulative regret over the first  $T_0$  rounds with  $T_0$ , and recalling that  $T_0 = \lceil T^{2/3} \rceil$ , yields the conclusion.  $\square$

### 3.2.2 A $T^{2/3}$ Lower Bound for the Smooth i.i.d. Environment

We prove that the  $\tilde{O}(T^{2/3})$  bound achieved by COLLECTING BANDIT is indeed optimal, up to logarithmic terms. Our lower bound consists in carefully embedding into our model a hard multiarmed bandit instance with  $K = \Theta(T^{1/3})$  arms, which entails a lower bound of order  $\Omega(\sqrt{KT}) = \Omega(T^{2/3})$ . This proof agenda involves various challenges: we want to embed a discrete construction of  $K$  independent actions into our continuous framework, where the utilities of different bids are correlated, while enforcing smoothness. Furthermore, the semi-transparent feedback is richer than the bandit one. We report here a proof sketch and refer the interested reader to Appendix A.2 for the missing details.

**Theorem 3.** *Consider the problem of repeated bidding in first-price auctions in a stochastic i.i.d.  $\sigma$ -smooth environment with semi-transparent feedback, for  $\sigma \in (0, 1/66]$ . Then, any learning algorithm  $\mathcal{A}$  satisfies, for  $T \geq 8$ ,*

$$R_T(\mathcal{A}) \geq \frac{3}{10^4} T^{2/3}.$$

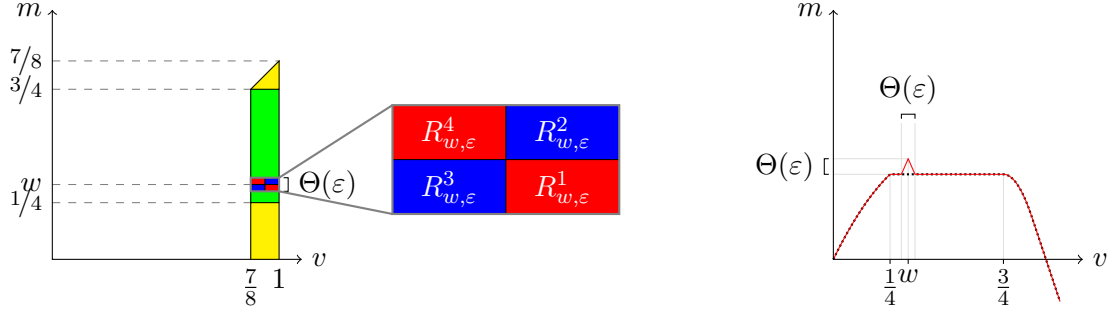


Figure 2: Left: The support of the base density  $f$  lies inside the yellow and green regions. The perturbation  $g_{w,\varepsilon}$  of  $f$  occurs inside the green region, where the four rectangles  $R_{w,\varepsilon}^1, \dots, R_{w,\varepsilon}^4$  (in red and blue) lie. Right: The corresponding qualitative plots of  $b \mapsto \mathbb{E}[\text{Util}_t(b)]$  (black, dotted) and  $p \mapsto \mathbb{E}^{w,\varepsilon}[\text{Util}_t(b)]$  (red, solid).

*Proof sketch.* Define, for all  $v, m \in [0, 1]$ , the density

$$f(v, m) = \mathbb{I}_{[\frac{7}{8}, 1]}(v) \left( \frac{1}{(v-m)^2} \mathbb{I}_{[\frac{1}{4}, v-\frac{1}{8}]}(m) + \frac{4}{v-1/4} \mathbb{I}_{[0, \frac{1}{4}]}(m) \right)^\S.$$

Let  $\mathbb{P}^0$  be a probability measure such that  $(V, M), (V_1, M_1), \dots$  is a  $\mathbb{P}$ -i.i.d. sequence where each pair  $(V, M)$  has common probability density function  $f$ . Denoting by  $\mathbb{E}^0$  the expectation with respect to  $\mathbb{P}^0$ , we have, for any bid  $b \in [0, 1]$  and any time step  $t$

$$\begin{aligned} \mathbb{E}^0[\text{Util}_t(b)] &= b \left( \frac{1}{2} + (1-4b) \ln \frac{6}{5} \right) \mathbb{I}_{[0, \frac{1}{4}]}(b) + \frac{1}{8} \mathbb{I}_{[\frac{1}{4}, \frac{3}{4}]}(b) \\ &\quad - \left( 4b^2 - 6b + \frac{17}{8} \right) \mathbb{I}_{[\frac{3}{4}, \frac{7}{8}]}(b) + \left( \frac{15}{16} - b \right) \mathbb{I}_{[\frac{7}{8}, 1]}(b). \end{aligned}$$

This function grows with  $b$  on  $[0, 1/4]$ , has a plateau of maximizers  $[1/4, 3/4]$ , then decreases on  $(3/4, 1]$  (see Figure 2, right). We introduce the perturbation space  $\Xi$ :

$$\Xi = \left\{ (w, \varepsilon) \in [0, 1]^2 : w - \varepsilon \geq \frac{1}{4} \text{ and } w + \varepsilon \leq \frac{3}{4} \right\}$$

and define, for all  $(w, \varepsilon) \in \Xi$ , the four rectangles

$$\begin{aligned} R_{w,\varepsilon}^1 &= [15/16, 1] \times [w - \varepsilon, w), & R_{w,\varepsilon}^2 &= [15/16, 1] \times [w, w + \varepsilon), \\ R_{w,\varepsilon}^3 &= [7/8, 15/16] \times [w - \varepsilon, w), & R_{w,\varepsilon}^4 &= [7/8, 15/16] \times [w, w + \varepsilon). \end{aligned}$$

For all  $(w, \varepsilon) \in \Xi$ , we introduce the probability density function  $f_{w,\varepsilon}$  as follows  $f_{w,\varepsilon} = f + g_{w,\varepsilon}$ , where the perturbation  $g_{w,\varepsilon}$  is defined as follows

$$g_{w,\varepsilon}(v, m) = \frac{16}{9} \left( \mathbb{I}_{R_{w,\varepsilon}^1 \cup R_{w,\varepsilon}^4}(v, m) - \mathbb{I}_{R_{w,\varepsilon}^2 \cup R_{w,\varepsilon}^3}(v, m) \right).$$

We refer to the left plot in Figure 2 for a visualization of the support of the  $f_{w,\varepsilon}$ . For all  $(w, \varepsilon) \in \Xi$ , let  $\mathbb{P}^{w,\varepsilon}$  be a probability measure such that  $(V, M), (V_1, M_1), (V_2, M_2), \dots$  is a  $\mathbb{P}^{w,\varepsilon}$ -i.i.d. sequence where each pair  $(V, M)$  has common probability density function  $f_{w,\varepsilon}$ . Denoting by  $\mathbb{E}^{w,\varepsilon}$  the expectation with respect to  $\mathbb{P}^{w,\varepsilon}$ , we have, for any bid  $b \in [0, 1]$  and any  $t$

$$\mathbb{E}^{w,\varepsilon}[\text{Util}_t(b)] = \mathbb{E}^0[\text{Util}_t(b)] + \frac{\varepsilon}{144} \Lambda_{w,\varepsilon}(b)$$

<sup>\S</sup>Note, we use the notation  $\mathbb{I}_A(x)$  to denote the indicator function that has value 1 when  $x \in A$ , and 0 otherwise.

where  $\Lambda_{u,r}$  is the tent map centered at  $u$  with radius  $r$  defined as  $\Lambda_{u,r}(x) = \max\{1 - |x - u|/r, 0\}$ . In words, in a perturbed scenario  $\mathbb{P}^{w,\varepsilon}$  the expected utility is maximized at the peak of a spike centered at  $w$  with length and height  $\Theta(\varepsilon)$  perturbing the plateau area  $[1/4, 3/4]$  of maximum height (see Figure 2, right). Define, for all times  $t \in \mathbb{N}$ , the feedback function  $\psi_t: [0, 1] \rightarrow ([0, 1] \times \{\star\}) \cup (\{\star\} \times [0, 1])$ , as follows:

$$b \mapsto \begin{cases} (V_t, \star) & \text{if } b \geq M_t \\ (\star, M_t) & \text{if } b < M_t \end{cases}$$

and note that, in our semi-transparent feedback model, the feedback  $Z_t$  received after bidding  $B_t$  at time  $t$  is  $\psi_t(B_t)$ . Crucially, for each  $(w, \varepsilon) \in \Xi$  and each  $b \in [0, 1] \setminus [w - \varepsilon, w + \varepsilon]$ , the distribution of  $\psi_t(b)$  under  $\mathbb{P}^{w,\varepsilon}$  coincides with the distribution of  $\psi_t(b)$  under  $\mathbb{P}^0$ . In push-forward notation (for a refresher on push-forward measures, see Appendix A.1), it holds that

$$\mathbb{P}_{\psi_t(b)}^{w,\varepsilon} = \mathbb{P}_{\psi_t(b)}^0. \quad (1)$$

Now, let  $K \in \mathbb{N}$ ,  $\varepsilon = 1/(4K)$ ,  $w_k = 1/4 + (2k - 1)\varepsilon$  and  $\mathbb{P}^k = \mathbb{P}^{w_k,\varepsilon}$  (for each  $k \in [K]$ ). At a high level, we built a problem with two crucial properties: (i) we know in advance the region where the optimal bid belongs to (i.e., the interval  $[1/4, 3/4]$ ), but (ii) when the underlying scenario is determined by the probability measure  $\mathbb{P}^k$ , the learner has to detect inside this potentially optimal region where a spike of height (and length)  $\Theta(\varepsilon)$  occurs (to avoid suffering suffer  $\Omega(\varepsilon T)$  regret). This last task can be accomplished only by locating where the perturbation in the base probability measure occurs, which, given the feedback structure, can only be done by playing in the interval  $[w_k - \varepsilon, w_k + \varepsilon]$  if the underlying probability is  $\mathbb{P}^k$ , suffering instantaneous regret of order  $\varepsilon$  whenever the underlying probability is  $\mathbb{P}^j$ , with  $j \neq k$ . Given that we partitioned the potentially optimal region  $[1/4, 3/4]$  into  $\Theta(1/\varepsilon)$  disjoint intervals where these perturbations can occur, the feedback structure implies that each of these intervals deserves its dedicated exploration.

To better highlight this underlying structure, in Appendix A.2, we show that our problem is not easier than a simplified  $K$ -armed stochastic bandit problem, where the instances we consider are determined by the probability measures  $\mathbb{P}^1, \dots, \mathbb{P}^K$ . In this bandit problem, when the underlying probability measure is induced by some  $\mathbb{P}^k$ , the corresponding arm  $k$  has an expected reward  $\Theta(\varepsilon)$  larger than the others. Then, via an information-theoretic argument, we can show that any learner would need to spend at least order of  $1/\varepsilon^2$  rounds to explore each of the  $K$  arms (paying  $\Omega(\varepsilon)$  each time) or else, it would pay a regret  $\Omega(\varepsilon T)$ . Hence, the regret of any learner, in the worst case, is lower bounded by  $\Omega(\frac{K}{\varepsilon^2}\varepsilon + \varepsilon T) = \Omega(K^2 + \frac{T}{K})$  (recalling our choice of  $\varepsilon = 1/(4K)$ ). Picking  $K = \Theta(T^{1/3})$  yields a lower bound of order  $T^{2/3}$ . For all missing technical details, see Appendix A.2.  $\square$

### 3.3 I.I.D. – Transparent/Full Feedback

This section completes the study of the stochastic i.i.d. environment by determining the minimax regret when the learner has access to full or transparent feedback.

#### 3.3.1 A $\sqrt{T}$ Upper Bound for the i.i.d. Environment

While with semi-transparent feedback, the learning algorithm has to rely on dummy bids  $B_1 = \dots = B_{T_0} = 0$  to gather information about the distribution of the highest competing bids, with the transparent one, this information is collected for free at each bidding round. To use this extra information, we present a wrapper W.T.FPA (for a sequence of base learning algorithms for the transparent feedback model) whose purpose is restarting the learning process with a geometric step to update the set of candidate bids. We assume that each of the wrapped base algorithms  $\tilde{\mathcal{A}}_\tau$  can

take as input any finite subset  $\mathcal{X} \subset [0, 1]$  and returns bids in  $\mathcal{X}$ . Furthermore, for all  $T'$ , we let  $\mathcal{R}_{T'}(\tilde{\mathcal{A}}_\tau, \mathcal{X})$  be an upper bound on the regret over  $T'$  rounds of  $\tilde{\mathcal{A}}_\tau$  with input  $\mathcal{X}$  against the best fixed  $x \in \mathcal{X}$ . Formally, we require that for any two times  $T_0 < T_1$  such that  $T' = T_1 - T_0$ , the quantity  $\mathcal{R}_{T'}(\tilde{\mathcal{A}}_\tau, \mathcal{X})$  is an upper bound on  $\max_{x \in \mathcal{X}} \mathbb{E} [\sum_{t=T_0+1}^{T_1} \text{Util}_t(x) - \sum_{t=T_0+1}^{T_1} \text{Util}_t(B_t)]$ , where  $B_t \in \mathcal{X}$  is the sequence of prices played by  $\tilde{\mathcal{A}}_\tau$  (with input  $\mathcal{X}$ ) when started at round  $t = T_0 + 1$  and ran up to time  $T_1$ . Without loss of generality, we assume that  $T' \mapsto \mathcal{R}_{T'}(\tilde{\mathcal{A}}_\tau, \mathcal{X})$  is non-decreasing.

---

W.T.FPA (Wrapper for Transparent First-Price Auctions)

---

- 1: **input:** Base algorithms  $\tilde{\mathcal{A}}_1, \tilde{\mathcal{A}}_2, \dots$
  - 2: **initialization:**  $s \leftarrow 0$
  - 3: **for** each epoch  $\tau = 1, 2, \dots$  **do**
  - 4:    $\mathcal{X}_\tau \leftarrow \{0\} \cup \{M_1, \dots, M_s\}$  (with  $\mathcal{X}_1 = \{0\}$ )
  - 5:   Start  $\tilde{\mathcal{A}}_\tau$  with input  $\mathcal{X}_\tau$  and run it for  $t = s + 1, \dots, s + 2^{\tau-1}$
  - 6:   Update  $s \leftarrow s + 2^{\tau-1}$
- 

**Proposition 1.** *Consider the problem of repeated bidding in first-price auctions in a stochastic i.i.d. environment with transparent feedback. Then the regret of W.T.FPA run with base algorithms  $\tilde{\mathcal{A}}_1, \tilde{\mathcal{A}}_2, \dots$  satisfies*

$$R_T(\text{W.T.FPA}) \leq \sum_{\tau=2}^{\lceil \log_2(T+1) \rceil} \mathcal{R}_{2^{\tau-1}}(\tilde{\mathcal{A}}_\tau, \mathcal{X}_\tau) + 3 + 16(\sqrt{2} + 2)\sqrt{T \ln T}.$$

*Proof.* Fix an arbitrary epoch  $\tau \in \{2, \dots, \lceil \log_2(T+1) \rceil\}$ ; we want to bound the regret suffered there by W.T.FPA using Lemma 3. Using the notation of the lemma, let  $\mathcal{X} = \mathcal{X}_\tau$ ,  $K+1 = |\mathcal{X}|$ ,  $T_0 = \sum_{\tau'=1}^{\tau-1} 2^{\tau'-1} = 2^{\tau-1} - 1$  (the time passed from the beginning of epoch 1 up to and including the end of epoch  $\tau-1$ ),  $T_1 = \min\{T_0 + 2^{\tau-1}, T\}$  (the end of epoch  $\tau$ ), and let  $X_0 < X_1 < \dots < X_K$  be the distinct elements of  $\mathcal{X}$  in increasing order, where we note that  $X_0 = 0$ ,  $X_K \leq 1$ , and we set  $X_{K+1} = 2$ . Let  $\mathcal{H}_{T_0}$  be the history, including time  $T_0$ .

Applying Lemma 3 (together with the fact that the empirical frequency between any two consecutive values  $X_k$  and  $X_{k+1}$  is 0 by design), and exploiting the monotonicity of  $T' \mapsto \mathcal{R}_{T'}(\tilde{\mathcal{A}}_\tau, \mathcal{X}_\tau)$  for the last epoch (if  $T_0 + 2^{\tau-1} > T$ ), we obtain, for all  $b \in [0, 1]$  and  $\delta \in (0, 1)$ ,

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=T_0+1}^{T_1} \text{Util}_t(b) \right] &\leq \sum_{t=T_0+1}^{T_1} \mathbb{E} \left[ \text{Util}_t(X_{k_{\mathcal{X}}(b)}) \right] + 2^{\tau-1} \left( 8\sqrt{\frac{\ln(1/\delta)}{T_0}} + \delta \right) \\ &\leq \sum_{t=T_0+1}^{T_1} \mathbb{E} [\text{Util}_t(B_t)] + \mathcal{R}_{2^{\tau-1}}(\tilde{\mathcal{A}}_\tau, \mathcal{X}_\tau) + 2^{\tau-1} \left( 8\sqrt{\frac{\ln(1/\delta)}{2^{\tau-1}-1}} + \delta \right). \end{aligned}$$

Summing over epochs  $\tau \in \{2, \dots, \lceil \log_2(T+1) \rceil\}$ , upper bounding by 1 the regret incurred in the first epoch, and tuning  $\delta = 1/T$  yields the conclusion.  $\square$

Now we are only left to design appropriate base algorithms  $\tilde{\mathcal{A}}_1, \tilde{\mathcal{A}}_2, \dots$  for the transparent feedback to wrap W.T.FPA around.

**The Exp3.FPA algorithm.** To this end, we introduce the Exp3.FPA algorithm (designed to run with transparent feedback), which borrows ideas from online learning with feedback graphs

[Alon et al., 2017]. Similar algorithms for related settings have been previously proposed by Weed et al. [2016] and Feng et al. [2018]. For the familiar reader, note that our setting can be seen as an instance of online learning with strongly observable feedback graphs. In contrast to a black-box application of feedback-graph results, we shave off a logarithmic term (in the time horizon) by using a dedicated analysis. For any  $x \in [0, 1]$ , we denote by  $\delta_x$  the Dirac distribution centered at  $x$ .

---

Exp3.FPA

---

- 1: **input:** Finite  $\mathcal{X} \subset [0, 1]$  with maximum  $\bar{x}$ , exploration rate  $\gamma \in (0, 1)$
- 2: For all  $x \in \mathcal{X}$ , let  $w_1(x) \leftarrow 1$
- 3: **for** each round  $t = 1, 2, \dots$  **do**
- 4:   Post bid  $B_t \sim p_t \leftarrow (1 - \gamma) \frac{w_t}{\|w_t\|_1} + \gamma \delta_{\bar{x}}$
- 5:   For all  $x \in \mathcal{X}$ , define the reward estimate:

$$\hat{g}_t(x) \leftarrow (V_t - x) \mathbb{I}\{x \geq M_t\} \frac{\mathbb{I}\{M_t \leq B_t\}}{\sum_{y \geq M_t} p_t(y)}$$

- 6:   For all  $x \in \mathcal{X}$ , update the weight:

$$w_{t+1}(x) \leftarrow w_t(x) \exp(\gamma \hat{g}_t(x))$$


---

Note that the transparent feedback is sufficient to compute the reward estimates in Line 5.

**Proposition 2.** *Let  $\mathcal{X} \subset [0, 1]$  be a finite set,  $T \in \mathbb{N}$  a time horizon, and tune the exploration rate as  $\gamma = \sqrt{\ln(|\mathcal{X}|)/(e-1)T}$ . Then, the regret of Exp3.FPA against the best fixed bid in  $\mathcal{X}$  is*

$$\max_{x \in \mathcal{X}} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(x) - \sum_{t=1}^T \text{Util}_t(B_t) \right] \leq 2\sqrt{(e-1) \ln(|\mathcal{X}|)T}$$

*Proof.* Let  $\gamma > 0$ . Notice that, for each  $t \in \mathbb{N}$ , it holds that  $\sum_{y \geq M_t} p_t(y) \geq \gamma$ . It follows, for each  $x \in \mathcal{X}$  and  $t \in \mathbb{N}$ , that  $\gamma \hat{g}_t(x) \leq 1$ , and hence

$$\exp(\gamma \hat{g}_t(x)) \leq 1 + \gamma \hat{g}_t(x) + (e-2)\gamma^2(\hat{g}_t(x))^2.$$

Then, for each  $t \in \mathbb{N}$ ,

$$\frac{\|w_{t+1}\|_1}{\|w_t\|_1} = \sum_{x \in \mathcal{X}} \frac{w_t(x)}{\|w_t\|_1} \exp(\gamma \hat{g}_t(x)) \leq 1 + \sum_{x \in \mathcal{X}} \frac{w_t(x)}{\|w_t\|_1} \left( \gamma \hat{g}_t(x) + (e-2)\gamma^2(\hat{g}_t(x))^2 \right),$$

which implies

$$\ln \left( \frac{\|w_{t+1}\|_1}{\|w_t\|_1} \right) \leq \sum_{x \in \mathcal{X}} \frac{w_t(x)}{\|w_t\|_1} \left( \gamma \hat{g}_t(x) + (e-2)\gamma^2(\hat{g}_t(x))^2 \right) \leq \frac{\gamma}{1-\gamma} \sum_{x \in \mathcal{X}} p_t(x) \left( \hat{g}_t(x) + (e-2)\gamma(\hat{g}_t(x))^2 \right).$$

Now, for each  $t \in \mathbb{N}$ , let  $\mathcal{F}_t$  be the  $\sigma$ -algebra generated by  $p_t, V_t$  and  $M_t$  and denote by  $\mathbb{E}_t := \mathbb{E}[\cdot | \mathcal{F}_t]$ . First, notice that, for each  $t \in \mathbb{N}$  and each  $x \in \mathcal{X}$

$$\mathbb{E}_t[\hat{g}_t(x)] = \text{Util}_t(x), \quad \mathbb{E}_t \left[ \sum_{x \in \mathcal{X}} p_t(x) \hat{g}_t(x) \right] = \mathbb{E}[\text{Util}_t(B_t) | V_t, M_t],$$



and that

$$\mathbb{E}_t \left[ \sum_{x \in \mathcal{X}} p_t(x) (\hat{g}_t(x))^2 \right] \leq \mathbb{E}_t \left[ \sum_{x \in \mathcal{X}} p_t(x) \frac{\mathbb{I}\{x \geq M_t\} \mathbb{I}\{M_t \leq B_t\}}{(\sum_{y \geq M_t} p_t(y))^2} \right] = \mathbb{E}_t \left[ \sum_{x \in \mathcal{X}} p_t(x) \frac{\mathbb{I}\{x \geq M_t\}}{\sum_{y \geq M_t} p_t(y)} \right] = 1 .$$

It follows that, for each  $x \in \mathcal{X}$ ,

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(x) \right] - \ln(|\mathcal{X}|) &= \mathbb{E} \left[ \sum_{t=1}^T \hat{g}_t(x) \right] - \ln(|\mathcal{X}|) = \mathbb{E} \left[ \ln(w_{T+1}(x)) \right] - \ln(|\mathcal{X}|) \\ &\leq \mathbb{E} \left[ \ln \left( \frac{\|w_{T+1}\|_1}{\|w_1\|_1} \right) \right] = \sum_{t=1}^T \mathbb{E} \left[ \mathbb{E}_t \left[ \ln \left( \frac{\|w_{t+1}\|_1}{\|w_t\|_1} \right) \right] \right] \\ &\leq \frac{\gamma}{1-\gamma} \left( \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(B_t) \right] + (e-2)\gamma T \right) , \end{aligned}$$

which, after rearranging and upper bounding, yields

$$\mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(x) - \sum_{t=1}^T \text{Util}_t(B_t) \right] \leq \frac{\ln(|\mathcal{X}|)}{\gamma} + (e-1)\gamma T .$$

Selecting  $\gamma$  as in the statement of the theorem leads to the conclusion.  $\square$

Putting together Propositions 1 and 2 yields the desired rate.

**Theorem 4.** *Consider the problem of repeated bidding in first-price auctions in a stochastic i.i.d. environment with transparent feedback. Then there exists a learning algorithm  $\mathcal{A}$  such that*

$$R_T(\mathcal{A}) \leq 3 + 2(\sqrt{2} + 2)(\sqrt{2(e-1)} + 8)\sqrt{T \ln T} .$$

*Proof.* The statement of the theorem holds for W.T.FPA run with the base algorithm of each epoch  $\tau$  being *Exp3.FPA* tuned with  $\gamma = \gamma(\tau) = \sqrt{\ln(|\mathcal{X}_\tau|) / ((e-1)2^{\tau-1})}$ . Substituting the guarantees of Proposition 2 into those of Proposition 1 and recalling that  $|\mathcal{X}_\tau| \leq 2^{\tau-1}$  for each epoch  $\tau = 2, 3, \dots$ , yields the desired bound.  $\square$

### 3.3.2 A $\sqrt{T}$ Lower Bound for the i.i.d. Environment

We complement the positive result of Theorem 4 with a matching lower bound of order  $\sqrt{T}$ . The idea underlying our hard instance is to embed the well-known lower bound for prediction with (two) experts into our framework: we construct two smooth distributions that are “similar” but have two different optimal bids whose performance is separated so that no learner can identify the correct distribution without suffering less than  $\sqrt{T}$  regret.

**Theorem 5.** *Consider the problem of repeated bidding in first-price auctions in a stochastic i.i.d.  $\sigma$ -smooth environment with full feedback, for  $\sigma \in (0, 1/9]$ . Then, any learning algorithm  $\mathcal{A}$  satisfies*

$$R_T(\mathcal{A}) \geq \frac{1}{2048} \sqrt{T} .$$

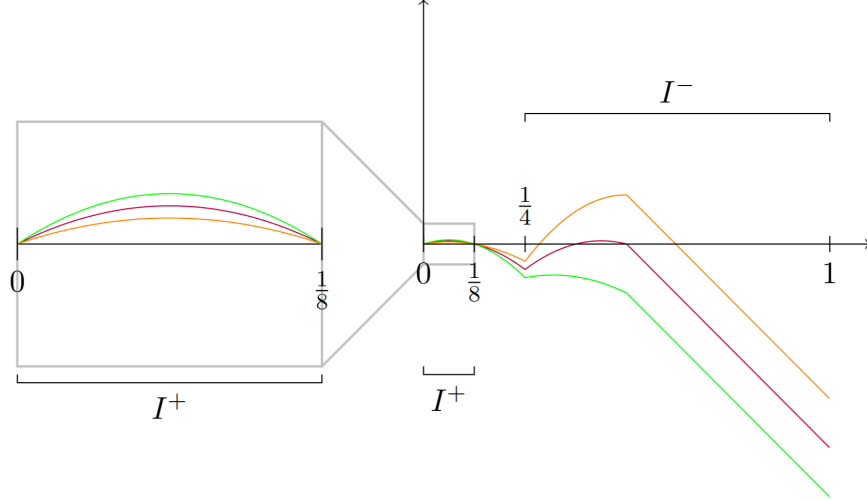


Figure 3: The expected utility function for three different distributions:  $\mathbb{P}^0$  in purple,  $\mathbb{P}^+$  in orange, and  $\mathbb{P}^-$  in green.

*Proof.* We prove the theorem by Yao's principle: we show that there exists a distribution over stochastic  $\sigma$ -smooth environments such that any deterministic learning algorithm  $\mathcal{A}$  suffers  $\Omega(\sqrt{T})$  regret against it, in expectation. We do that in two steps. First, for every  $\varepsilon \in (0, 1/2)$  we construct a pair of  $1/9$ -smooth distributions that are hard to discriminate for the learner. Then, we prove that, for the right choice of  $\varepsilon$ , any learner suffers the desired regret against a uniform mixture of them. For visualization, we refer to Figure 3.

As a tool for our construction, we introduce a baseline probability measure  $\mathbb{P}^0$ , such that the sequence  $(V, M), (V_1, M_1), (V_2, M_2), \dots$  is  $\mathbb{P}^0$ -i.i.d., and  $(V, M)$  has distribution  $\mathbb{P}_{(V, M)}^0$  (for a refresher on push-forward measures, see Appendix A.1) whose pdf is

$$f^0(v, m) = 8(\mathbb{I}_{Q_+}(v, m) + \mathbb{I}_{Q_-}(v, m)),$$

where  $Q_+ = (0, 1/4) \times (0, 1/4)$  and  $Q_- = (3/4, 1) \times (1/4, 1/2)$ . A convenient way to visualize this distribution is to draw a uniform random variable  $U_t$  in the square  $Q_+$  and then toss an unbiased coin. If the coin yields heads, then  $(V_t, M_t)$  is equal to  $U_t$ , otherwise  $(V_t, M_t)$  coincides with  $U_t$  translated by  $(3/4, 1/4)$ . With some simple computation, it is possible to explicitly compute the expected utility of posting any bid  $b \in [0, 1]$  when  $(V_t, M_t)$  is drawn following the distribution  $\mathbb{P}^0$  (and expectation  $\mathbb{E}^0$ ):

$$\mathbb{E}^0[\text{Util}_t(b)] = \begin{cases} \frac{b}{4}(1 - 8b) & \text{if } b \in [0, 1/4) \\ -\frac{1}{8}(16b^2 - 14b + 3) & \text{if } b \in [1/4, 1/2) \\ \frac{1}{2}(1 - 2b) & \text{if } b \in [1/2, 1] \end{cases}$$

The function  $\mathbb{E}^0[\text{Util}_t(b)]$  has two global maxima in  $[0, 1]$ , of value  $1/128$ , attained in  $1/16$  and  $7/16$  (see purple line in Figure 3).

For any  $\varepsilon \in (0, 1/2)$ , we also define two additional (perturbed) probability measures  $\mathbb{P}^{\pm\varepsilon}$ , such that the sequence  $(V, M), (V_1, M_1), \dots$  is  $\mathbb{P}^{\pm\varepsilon}$ -i.i.d. and the distribution  $\mathbb{P}_{(V, M)}^{\pm\varepsilon}$  of  $(V, M)$  has density:

$$f^{\pm\varepsilon}(v, m) = 8(1 \pm \varepsilon)\mathbb{I}_{Q_+}(v, m) + 8(1 \mp \varepsilon)\mathbb{I}_{Q_-}(v, m).$$

Note,  $\|f^{\pm\varepsilon}\|_\infty < 9$ , while  $\|f^0\|_\infty = 8$ , therefore all the distributions considered in this proof are  $1/9$ -smooth. To visualize these new perturbed distributions, recall the construction of  $\mathbb{P}_{(V, M)}^0$  using

the coin toss and the uniform random variable  $U$ : in this case, the coin is biased, and the probability of tails is  $(1 \pm \varepsilon)/2$ . It is possible to explicitly compute the expected utility under these perturbed distributions for any bid  $b \in [0, 1]$ :  $\mathbb{E}^{\pm \varepsilon}[\text{Util}_t(b)]$  is equal to

$$\begin{cases} \frac{b}{4}(1 - 8b) \pm \varepsilon \frac{b}{4}(1 - 8b) & \text{if } b \in [0, \frac{1}{4}) \\ -\frac{1}{8}(16b^2 - 14b + 3) \pm \frac{\varepsilon}{4}(8b^2 - 11b + 2) & \text{if } b \in [\frac{1}{4}, \frac{1}{2}) \\ \frac{1}{2}(1 - 2b \mp \frac{3}{4}\varepsilon) & \text{if } b \in [\frac{1}{2}, 1] \end{cases} \quad (2)$$

For visualization, we refer to Figure 3 (bottom). The crucial property of the distributions we constructed is that the instantaneous regret of not playing in the “correct” region is  $\Omega(\varepsilon)$ ; formally we have the following result. For the sake of readability, we postpone the proof of this claim to Appendix A.3.

**Claim 1.** *There exists two disjoint intervals  $I_+$  and  $I_-$  in  $[0, 1]$  such that, for any  $\varepsilon \in (0, 1/2)$  and any time  $t$ , the following hold:*

$$\max_{x \in [0, 1]} \mathbb{E}^{\pm \varepsilon}[\text{Util}_t(x)] \geq \mathbb{E}^{\pm \varepsilon}[\text{Util}_t(b)] + \frac{1}{128}\varepsilon, \text{ for all } b \notin I_{\pm}$$

Since the two distributions are “ $\varepsilon$ -close<sup>¶</sup>”, any learner needs at least  $1/\varepsilon^2$  rounds to discriminate which ones of the two distributions it is actually facing, paying each error with an instantaneous regret of  $\Omega(\varepsilon)$  (Claim 1). All in all, any learner suffers a regret that is  $\Omega(\varepsilon \cdot \frac{1}{\varepsilon^2} + \varepsilon T)$ , which is of the desired  $\Omega(\sqrt{T})$  order for the right choice of  $\varepsilon \approx T^{-1/2}$ .

As the last step of the proof, we formalize the above argument. Fix  $\varepsilon = 1/(4\sqrt{T})$  and rename  $\mathbb{P}^{+\varepsilon} = \mathbb{P}^1$  and  $\mathbb{P}^{-\varepsilon} = \mathbb{P}^2$ . Similarly, denote with  $I_+$  and  $I_-$  the two intervals  $I_+$  and  $I_-$  as in the statement of Claim 1. For each  $j \in \{0, 1, 2\}$ , consider the run of  $\mathcal{A}$  against the stochastic environment which draws  $(V_1, M_1), (V_2, M_2), \dots$  i.i.d. from  $\mathbb{P}^j$ . Let  $N_1$  be the random variable that counts the number of times that algorithm  $\mathcal{A}$  posts a bid in  $I_+$ . Similarly,  $N_2$  counts the number of times that it posts a bid in  $I_-$ . For  $i = 1, 2$ , we have the following crucial relation between the expected value of  $N_i$  under  $\mathbb{P}^i$ . Note, the results hold because the two distributions are so similar that the deterministic algorithm  $\mathcal{A}$  bids in the wrong region a constant fraction of the time steps. For the formal proof of we refer the reader to Appendix A.3.

**Claim 2.** *The following inequality holds:  $\frac{1}{2} \sum_{i=1,2} \mathbb{E}^i [N_i] \leq \frac{3}{4}T$ .*

We finally have all the ingredients to conclude the proof. Consider an environment that selects uniformly at random either  $\mathbb{P}^1$  or  $\mathbb{P}^2$  and then draws the  $(V_t, M_t)$  i.i.d. following it. We prove that the algorithm  $\mathcal{A}$  suffers linear regret against this randomized environment and, by a simple averaging argument, against at least one of them. Specifically, if  $b_i^*$  is the optimal bid in the scenario determined by  $\mathbb{P}^i$ , for  $i \in \{1, 2\}$ , we have

$$\begin{aligned} R_T(\mathcal{A}) &\geq \frac{1}{2} \sum_{i=1,2} \mathbb{E}^i \left[ \sum_{t=1}^T \text{Util}_t(b_i^*) - \sum_{t=1}^T \text{Util}_t(B_t) \right] \\ &\stackrel{(*)}{\geq} \frac{1}{1024\sqrt{T}} \sum_{i=1,2} \mathbb{E}^i [T - N_i] \stackrel{(\circ)}{\geq} \frac{1}{512\sqrt{T}} \left( T - \frac{3}{4}T \right) = \frac{\sqrt{T}}{2048} \end{aligned}$$

where  $(*)$  follows by Claim 1 and choice of  $\varepsilon$ , and  $(\circ)$  by Claim 2. □

<sup>¶</sup>In Appendix A.3 we formally prove that their total variation is at most  $\Theta(\varepsilon)$ .

## 4 The Adversarial Setting

In this section we complete the perspective on repeated bidding in first-price auction by investigating the adversarial environment. In particular, we consider two models: the standard one, where the sequence  $(V_1, M_1), (V_2, M_2), \dots$  is chosen upfront in a deterministic oblivious way, and the smooth environment, where the sequence  $(V_1, M_1), (V_2, M_2), \dots$  is some  $\sigma$ -smooth stochastic process. In Section 4.1 we construct an algorithm achieving  $T^{2/3}$  regret in the bandit feedback model under the smoothness assumption; this result, together with the lower bound of the same order for the semi-transparent feedback (Theorem 3) settles the problem for these two feedback regimes. Then, in Section 4.2 we provide another upper bound, namely an algorithm achieving  $\sqrt{T}$  regret in the transparent feedback model under the smoothness assumption; this result, together with the lower bound of the same order for the semi-transparent feedback (Theorem 5) settles the problem for these two feedback regimes. Finally, in Section 4.3 we provide a lower bound proving that the non-smooth adversarial environment is too hard to learn, even when the learner has access to full feedback.

### 4.1 Smooth – Bandit Feedback

The smoothness assumption regularizes the objective function: if  $(V_t, M_t)$  is smooth, then the expected utility is Lipschitz.

**Lemma 4** (Lipschitzness). *Let  $(V_t, M_t)$  be a  $\sigma$ -smooth random variable in  $[0, 1]^2$ . Then the induced expected utility function  $\mathbb{E}[\text{Util}_t(\cdot)]$  is  $2/\sigma$ -Lipschitz in  $[0, 1]$ :*

$$|\mathbb{E}[\text{Util}_t(y) - \text{Util}_t(x)]| \leq \frac{2}{\sigma}|y - x|, \quad \forall x, y \in [0, 1]. \quad (3)$$

*Proof.* Let  $x > y$  be any two bids in  $[0, 1]$ , we have:

$$\begin{aligned} |\mathbb{E}[\text{Util}_t(x) - \text{Util}_t(y)]| &= |\mathbb{E}[(V_t - x)\mathbb{I}\{M_t \leq x\} - (V_t - y)\mathbb{I}\{M_t \leq y\}]| \\ &= |\mathbb{E}[(V_t - x)\mathbb{I}\{y < M_t \leq x\} + (x - y)\mathbb{I}\{M_t \leq y\}]| \\ &\leq \mathbb{P}[M_t \in [x, y]] + (x - y) \leq \frac{2}{\sigma}(x - y). \quad \square \end{aligned}$$

Interestingly, we only need the marginal distribution of  $M_t$  to be  $\sigma$ -smooth for the previous lemma to hold. This Lipschitzness property has the immediate corollary that any fine enough discretization of  $[0, 1]$  contains a bid whose utility is close to the optimal one.

**Lemma 5** (Discretization Lemma). *Let  $\mathcal{X}$  be any finite grid of bids in  $[0, 1]$ , and let  $\delta(\mathcal{X})$  be the largest distance of a point in  $[0, 1]$  to  $\mathcal{X}$  (i.e.,  $\delta(\mathcal{X}) = \max_{p \in [0, 1]} \min_{x \in \mathcal{X}} |p - x|$ ), then if each pair of random variables  $(V_1, M_1), \dots, (V_T, M_T)$  is  $\sigma$ -smooth, we have the following:*

$$\sup_{b \in [0, 1]} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(b) \right] - \max_{x \in \mathcal{X}} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(x) \right] \leq 3 \frac{\delta(\mathcal{X})}{\sigma} T.$$

*Proof.* Fix any such sequence and let  $b^*$  a fixed bid such that

$$\sup_{b \in [0, 1]} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(b) \right] \leq \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(b^*) \right] + \frac{\delta(\mathcal{X})}{\sigma} T. \quad (4)$$

If  $b^*$  is in  $\mathcal{X}$  there is nothing to prove, otherwise there exists  $x^* \in \mathcal{X}$  such that  $|b^* - x^*| \leq \delta(\mathcal{X})$  (by definition of  $\delta(\mathcal{X})$ ). It holds that

$$\sum_{t=1}^T \mathbb{E}[\text{Util}_t(b^*) - \text{Util}_t(x^*)] \stackrel{(L)}{\leq} \sum_{t=1}^T \frac{2}{\sigma} |b^* - x^*| \leq 2 \frac{\delta(\mathcal{X})}{\sigma} T.$$

where (L) follows by Lipschitzness and Lemma 4. The right-hand side with Equation (4) concludes the proof of the lemma.  $\square$

We can combine the above discretization lemma with any (optimal) bandits algorithm to get the desired bound on the regret. For details, we refer to the pseudocode of DISCRETIZED BANDIT.

---

DISCRETIZED BANDIT

---

- 1: **input:** Time horizon  $T$ , bandit algorithm  $\tilde{\mathcal{A}}$  for gains in  $[-1, 1]$ , grid of  $K$  bids  $\mathcal{X}$
  - 2: Initialize  $\tilde{\mathcal{A}}$  on  $K$  actions, one for each  $x \in \mathcal{X}$ , time horizon  $T$
  - 3: **for** each round  $t = 1, 2, \dots, T$  **do**
  - 4:     Receive from  $\tilde{\mathcal{A}}$  the bid  $B_t \in \mathcal{X}$
  - 5:     Post bid  $B_t$  and observe feedback  $Z_t$
  - 6:     Reconstruct  $\text{Util}_t(B_t)$  from  $Z_t$  and feed it to  $\tilde{\mathcal{A}}$
- 

**Theorem 6.** *Consider the problem of repeated bidding in first-price auctions in an adversarial  $\sigma$ -smooth environment with bandit feedback. Then there exists a learning algorithm  $\mathcal{A}$  such that*

$$R_T(\mathcal{A}) \leq \frac{29}{\sigma} T^{2/3}.$$

*Proof.* We prove that algorithm DISCRETIZED BANDIT with the right choice of learning algorithm  $\tilde{\mathcal{A}}$  and grid of bids  $\mathcal{X}$  achieves the desired bound on the regret. As learning algorithm  $\tilde{\mathcal{A}}$  we use (a rescaled version of) the Poly INF algorithm [Audibert and Bubeck, 2010]: since Poly INF is designed to run with gains in  $[0, 1]$  while the utilities we observe are in  $[-1, 1]$ , we first apply the reward transformation  $x \mapsto \frac{x+1}{2}$  to the observed utilities. This transformation costs a multiplicative factor of 2 in the regret guarantees of Poly INF.

The analysis builds on the discretization result in Lemma 5, by choosing as  $\mathcal{X}$  the uniform grid of  $\lceil T^{2/3} \rceil + 1$  equally spaced bids on  $[0, 1]$  (note,  $\delta(\mathcal{X})$  becomes  $T^{-1/3}$ ). Fix any  $\sigma$ -smooth environment  $\mathcal{S}$ , by Lemma 5, the following chain of inequalities holds:

$$\begin{aligned} \max_{b \in [0,1]} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(b) \right] &\leq \max_{x \in \mathcal{X}} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(x) \right] + \frac{6}{\sigma} T^{2/3} \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(B_t) \right] + \frac{6}{\sigma} T^{2/3} + 23T^{2/3} \leq \frac{29}{\sigma} T^{2/3}. \end{aligned}$$

The second inequality follows from the guarantees of (the rescaled version of) Poly INF [Audibert and Bubeck, 2010, Theorem 11].  $\square$

## 4.2 Smooth – Transparent Feedback

For transparent feedback, we combine two tools: the adversarial discretization result (Lemma 5) and the algorithm Exp3.FPA for learning with transparent feedback on a finite grid. Note, using any other  $\sqrt{KT}$  black box learning algorithm (like in the previous section for bandits) would yield a suboptimal regret bound of  $T^{2/3}$ .

**Theorem 7.** *Consider the problem of repeated bidding in first-price auctions in an adversarial  $\sigma$ -smooth environment with transparent feedback. Then there exists a learning algorithm  $\mathcal{A}$  such that*

$$R_T(\mathcal{A}) \leq 6 \left( \frac{1}{\sigma} + \sqrt{\ln T} \right) \sqrt{T}.$$

*Proof.* Consider algorithm Exp3.FPA on the uniform grid  $\mathcal{X}$  of  $\lceil \sqrt{T} \rceil + 1$  bids, with  $\delta(\mathcal{X}) \leq \sqrt{T}$ . Fix any  $\sigma$ -smooth environment  $\mathcal{S}$ , Lemma 5 implies the following:

$$\begin{aligned} \max_{b \in [0,1]} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(b) \right] &\leq \max_{x \in \mathcal{X}} \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(x) \right] + \frac{6}{\sigma} \sqrt{T} \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \text{Util}_t(B_t) \right] + 6 \left( \frac{1}{\sigma} + \sqrt{\ln T} \right) \sqrt{T}, \end{aligned}$$

where the second inequality follows from Proposition 2.  $\square$

### 4.3 The (Non-Smooth) Adversarial Model

The positive results provided in the previous sections hold under either one of two conditions: the environment is stochastic and the learner has at least the semi-transparent feedback (Theorem 1 says that bandit feedback is not enough) or the environment uses smooth distributions. These settings allow the learner to compute a discrete class of representative bids efficiently. In this section, we formally argue that learning is impossible if any of these assumptions is dropped. Specifically, the standard adversarial environment that generates the sequence without any smoothness constraint is too strong. In particular, we construct a randomized sequence  $(V_1, M_1), (V_2, M_2), \dots$  that induces any learner to suffer at least linear regret. This construction shares some similarities with the lower bound construction in Theorem 1, the main difference being that the best bid  $b^*$  is randomized and hidden in such a way that even a learner having access to full feedback cannot pin-point it.

**Theorem 8.** *Consider the problem of repeated bidding in first-price auctions in an adversarial environment with full feedback. Then, any learning algorithm  $\mathcal{A}$  satisfies  $R_T(\mathcal{A}) \geq T/24$ .*

*Proof.* We prove the result via Yao’s principle, showing that there exists a randomized environment  $\mathcal{S}$  such that any deterministic learning algorithm suffers  $T/24$  regret against it. The random sequence posted by  $\mathcal{S}$  is based on two randomized auxiliary sequences  $L_1, L_2, \dots$  and  $U_1, U_2, \dots$  defined as follows. They are initiated to  $L_0 = 1/2, U_0 = 2/3$ . They then evolve recursively as follows:

$$\begin{cases} L_t = L_{t-1} + \frac{2}{3}\Delta_{t-1} \text{ and } U_t = U_{t-1}, \text{ with probability } \frac{1}{2}, \\ U_t = U_{t-1} - \frac{2}{3}\Delta_{t-1} \text{ and } L_t = L_{t-1}, \text{ with probability } \frac{1}{2}, \end{cases}$$

where  $\Delta_{t-1} = U_{t-1} - L_{t-1}$ . For each realized sequence of the  $(L_t, U_t)$  pairs, the actual sequence of the  $(M_t, V_t)$  selected by  $\mathcal{S}$  is constructed as follows. At each time step  $t$ , the environment selects  $(M_t, V_t) = (L_t, 1)$  or  $(U_t, 0)$ , uniformly at random; so that the distribution is characterized by two levels of independent randomness: the auxiliary sequence of shrinking intervals and the choice between  $(L_t, 1)$  and  $(U_t, 0)$ .

We move our attention to the expected performance of the best fixed bid in hindsight. For each realization of the random auxiliary sequence, there exists a bid  $B^*$  such that (i) it wins all the auctions  $(V_t, M_t)$  of the form  $(L_t, 1)$  (which we may call “good auctions” because they bring positive utility when won) and (ii) it loses all the auctions  $(V_t, M_t)$  of the form  $(U_t, 0)$  (called “bad auctions” because they bring negative utility). Thus its expected utility at each time step is at least  $1/6$ : with probability  $1/2$  the environment selects a good auction, which induces a utility of  $(1 - L_t) \geq 1/3$ . All in all, the optimal bid achieves an expected utility of at least  $T/6$ .

Consider now the performance of any deterministic algorithm  $\mathcal{A}$ : for any fixed time  $t > 1$  and possible realization of the past observations, the learner posts some deterministic bid  $B_t$ . If  $B_t < L_{t-1}$ , then it gets 0 utility, so we only consider the following cases:

- If  $B_t \in [L_{t-1}, L_{t-1} + \frac{1}{3}\Delta_{t-1})$ , then the bidder gets the item with probability  $1/4$  ( $L_t = L_{t-1}$ ,  $V_t$  is set to 1 and  $M_t = L_t$ ) with an expected utility of  $(1-L_t)/4 \leq 1/8$ .
- If  $B_t \in [L_{t-1} + \frac{1}{3}\Delta_{t-1}, L_{t-1} + \frac{2}{3}\Delta_{t-1})$ , the bidder gets the item with probability  $1/2$  (when  $L_t = L_{t-1}$  and  $U_t = U_{t-1} - \frac{2}{3}\Delta_{t-1}$ ) for an expected utility of  $\frac{1}{4}(1 - 2L_{t-1} - \frac{1}{3}\Delta_{t-1}) \leq 0$
- If  $B_t \in [L_{t-1} + \frac{2}{3}\Delta_{t-1}, U_{t-1})$ , the bidder gets the item with probability  $3/4$  (when  $L_t = L_{t-1}$  and when  $U_t = U_{t-1}$ ,  $V_t = 1$  and  $M_t = L_t$ ) for an expected utility of  $\frac{1}{4}(1 - L_{t-1}) - \frac{1}{4}(L_{t-1} + \frac{1}{3}\Delta_{t-1}) + \frac{1}{4}(1 - L_{t-1} - \frac{2}{3}\Delta_{t-1}) \leq \frac{1}{8}$
- If  $B_t \geq U_{t-1}$ , the bidder always gets the item with a negative expected utility.

All in all, the expected utility of any deterministic algorithm is at most  $T/8$ . If we compare this quantity with the lower bound on the expected utility of the best bid in hindsight, we get the desired result:  $\mathbb{E}[R_T(\mathcal{A}, \mathcal{S})] \geq T/6 - T/8 = T/24$ .  $\square$

A final observation: the main ingredient in the proof is the elaborate auxiliary sequence. To construct it, we only needed the non-smoothness of  $M_t$ , while we may have chosen the valuations  $V_t$  to be smooth, say uniformly in  $[0, 1/4]$  for the bad auctions and in  $[3/4, 1]$  for the good ones.

## 5 Conclusion

Motivated by the recent shift from second to first-price auctions in online advertising markets, this paper comprehensively analyzes the online learning problem of repeated bidding in first-price auctions under the realistic assumption that the bidder does not know its valuation before bidding. We characterize the minimax regret achievable for different levels of transparency in the auction format and different data generation models, considering both the stochastic i.i.d. and the standard adversarial model, while also considering smoothness. Although our regret rates are tight in their dependence on the time horizon  $T$ , a natural open problem is studying their minimax dependence on the smoothness parameter  $\sigma$ . This paper belongs to the long line of research that studies economic problems from the online learning perspective; an intriguing open problem consists in offering a unified framework to characterize in a satisfying way all these games with partial feedback, similar to what has been done for partial monitoring and feedback graphs.

## Acknowledgment

NCB, RC, FF, and SL are partially supported by the FAIR (Future Artificial Intelligence Research) project, funded by theNextGenerationEU program within the PNRR-PE-AI scheme (M4C2, investment 1.3, line on Artificial Intelligence). NCB and RC are also partially supported by the MUR PRIN grant 2022EKNE5K (Learning in Markets and Society) and by the EU Horizon CL4-2022-HUMAN-02 RIA under grant agreement 101120237, project ELIAS (European Lighthouse of AI for Sustainability). RC also acknowledges the financial support of the Italian Institute of Technology during the writing of this paper. FF and SL are also partially supported by ERC Advanced Grant 788893 AMDROMA and PNRR MUR project IR0000013-SoBigData.it.

TC gratefully acknowledges the support of the University of Ottawa through grant GR002837 (Start-Up Funds) and that of the Natural Sciences and Engineering Research Council of Canada (NSERC) through grants RGPIN-2023-03688 (Discovery Grants Program) and DGEGR-2023-00208 (Discovery Grants Program, DGEGR - Discovery Launch Supplement)

## References

- Juliette Achddou, Olivier Cappé, and Aurélien Garivier. Fast rate learning in stochastic first price bidding. In *ACML*, volume 157 of *Proceedings of Machine Learning Research*, pages 1754–1769. PMLR, 2021.
- Noga Alon, Nicolò Cesa-Bianchi, Claudio Gentile, Shie Mannor, Yishay Mansour, and Ohad Shamir. Nonstochastic multi-armed bandits with graph-structured feedback. *SIAM J. Comput.*, 46(6): 1785–1826, 2017. doi: 10.1137/140989455.
- Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, 2009.
- Jean-Yves Audibert and Sébastien Bubeck. Regret bounds and minimax policies under partial monitoring. *J. Mach. Learn. Res.*, 11:2785–2836, 2010.
- Ashwinkumar Badanidiyuru, Zhe Feng, and Guru Guruganesh. Learning to bid in contextual first price auctions. In *WWW*, pages 3489–3497. ACM, 2023.
- Maria-Florina Balcan, Travis Dick, and Ellen Vitercik. Dispersion for data-driven algorithm design, online learning, and private optimization. In *FOCS*, pages 603–614. IEEE Computer Society, 2018. doi: 10.1109/FOCS.2018.00064.
- Santiago R. Balseiro, Negin Golrezaei, Mohammad Mahdian, Vahab S. Mirrokni, and Jon Schneider. Contextual bandits with cross-learning. *NeurIPS*, 2019.
- Gábor Bartók, Dean P. Foster, Dávid Pál, Alexander Rakhlin, and Csaba Szepesvári. Partial monitoring - classification, regret bounds, and algorithms. *Math. Oper. Res.*, 39(4):967–997, 2014. doi: 10.1287/moor.2014.0663.
- Richard F. Bass. *Real analysis for graduate students*. Createspace Ind Pub, 2013.
- Dirk Bergemann and Johannes Hörner. Should first-price auctions be transparent? *American Economic Journal: Microeconomics*, 10(3):177–218, 2018.
- Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. No-regret learning in bilateral trade via global budget balance. In *STOC*. ACM, 2024.
- Jason Bigler. Rolling out first price auctions to Google Ad Manager partners. <https://www.blog.google/products/admanager/rolling-out-first-price-auctions-google-ad-manager-partners/>, 2019. Accessed April 7, 2023.
- Adam Block, Yuval Dagan, Noah Golowich, and Alexander Rakhlin. Smoothed online learning is as easy as statistical learning. In *COLT*, volume 178 of *Proceedings of Machine Learning Research*, pages 1716–1786. PMLR, 2022.
- Nataša Bolić, Tommaso Cesari, and Roberto Colomboni. An online learning theory of brokerage. *The 23rd International Conference on Autonomous Agents and Multi-Agent Systems*, 2024.
- Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. A regret analysis of bilateral trade. In *EC*, pages 289–309. ACM, 2021. doi: 10.1145/3465456.3467645.



- Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. Repeated bilateral trade against a smoothed adversary. In *COLT*, volume 195 of *Proceedings of Machine Learning Research*, pages 1095–1130. PMLR, 2023.
- Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. The role of transparency in repeated first-price auctions with unknown valuations. In *STOC*. ACM, 2024a.
- Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. Regret analysis of bilateral trade with a smoothed adversary. *hal preprint hal-04383576*, 2024b.
- Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. Bilateral trade: A regret minimization perspective. *Mathematics of Operations Research*, 49(1): 171–203, 2024c. doi: 10.1287/moor.2023.1351.
- Xiaotie Deng, Xinyan Hu, Tao Lin, and Weiqiang Zheng. Nash convergence of mean-based learning algorithms in first price auctions. In *WWW*. ACM, 2022.
- Nishanth Dikkala and Éva Tardos. Can credit increase revenue? In *WINE*, volume 8289 of *Lecture Notes in Computer Science*, pages 121–133. Springer, 2013.
- Paul Duetting, Guru Guruganesh, Jon Schneider, and Joshua Ruizhi Wang. Optimal no-regret learning for one-sided lipschitz functions. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 8836–8850. PMLR, 2023.
- Naveen Durvasula, Nika Haghtalab, and Manolis Zampetakis. Smoothed analysis of online non-parametric auctions. In *EC*, pages 540–560. ACM, 2023.
- Michal Feldman, Brendan Lucier, and Noam Nisan. Correlated and coarse equilibria of single-item auctions. In *WINE*, volume 10123 of *Lecture Notes in Computer Science*, pages 131–144. Springer, 2016. doi: 10.1007/978-3-662-54110-4\_10.
- Zhe Feng, Chara Podimata, and Vasilis Syrgkanis. Learning to bid without knowing your value. In *EC*, pages 505–522. ACM, 2018.
- Zhe Feng, Guru Guruganesh, Christopher Liaw, Aranyak Mehta, and Abhishek Sethi. Convergence analysis of no-regret bidding algorithms in repeated auctions. In *AAAI*, pages 5399–5406. AAAI Press, 2021. doi: 10.1609/aaai.v35i6.16680.
- Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis of online and differentially private learning. In *NeurIPS*, 2020.
- Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis with adaptive adversaries. In *FOCS*, pages 942–953. IEEE, 2021.
- Nika Haghtalab, Yanjun Han, Abhishek Shetty, and Kunhe Yang. Oracle-efficient online learning for smoothed adversaries. In *NeurIPS*, 2022.
- Yanjun Han, Zhengyuan Zhou, Aaron Flores, Erik Ordentlich, and Tsachy Weissman. Learning to bid optimally and efficiently in adversarial first-price auctions. *arXiv preprint arXiv:2007.04568*, 2020a.
- Yanjun Han, Zhengyuan Zhou, and Tsachy Weissman. Optimal no-regret learning in repeated first-price auctions. *arXiv preprint arXiv:2003.09795*, 2020b.

- Sampath Kannan, Jamie H Morgenstern, Aaron Roth, Bo Waggoner, and Zhiwei Steven Wu. A smoothed analysis of the greedy algorithm for the linear contextual bandit problem. *Advances in neural information processing systems*, 31, 2018.
- Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. Bandits and experts in metric spaces. *J. ACM*, 66(4):30:1–30:77, 2019. doi: 10.1145/3299873.
- Yoav Kolumbus and Noam Nisan. Auctions between regret-minimizing agents. In *WWW*, pages 100–111. ACM, 2022. doi: 10.1145/3485447.3512055.
- Tor Lattimore. Minimax regret for partial monitoring: Infinite outcomes and rustichini’s regret. In *COLT*, volume 178 of *Proceedings of Machine Learning Research*, pages 1547–1575. PMLR, 2022.
- Michael Mitzenmacher and Eli Upfal. *Probability and Computing: Randomized Algorithms and Probabilistic Analysis, Second Edition*. Cambridge University Press, 2017. doi: 10.1017/CBO9780511813603.
- Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Stochastic, constrained, and smoothed adversaries. In *NIPS*, 2011.
- Dravyansh Sharma, Maria-Florina Balcan, and Travis Dick. Learning piecewise lipschitz functions in changing environments. In *AISTATS*, volume 108 of *Proceedings of Machine Learning Research*, pages 3567–3577. PMLR, 2020.
- Aleksandrs Slivkins. Introduction to multi-armed bandits. *Found. Trends Mach. Learn.*, 12(1-2): 1–286, 2019. doi: 10.1561/22000000068.
- Sarah Sluis. Big changes coming to auctions, as exchanges roll the dice on first-price. <https://adexchanger.com/platforms/big-changes-coming-auctions-exchanges-roll-dice-first-price/>, 2017. Accessed July 3, 2023.
- Daniel A Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM (JACM)*, 51(3):385–463, 2004. doi: 10.1145/990308.990310.
- Jonathan Weed, Vianney Perchet, and Philippe Rigollet. Online learning in repeated auctions. In *COLT*, volume 49 of *JMLR Workshop and Conference Proceedings*, pages 1562–1583. JMLR.org, 2016.
- Matt Wong. Moving AdSense to a first-price auction. <https://blog.google/products/ads-commerce/our-move-to-a-first-price-auction/>, 2021. Accessed July 6, 2023.
- Wei Zhang, Brendan Kitts, Yanjun Han, Zhengyuan Zhou, Tingyu Mao, Hao He, Shengjun Pan, Aaron Flores, San Gultekin, and Tsachy Weissman. MEOW: A space-efficient nonparametric bid shading algorithm. In *KDD*. ACM, 2021.
- Wei Zhang, Yanjun Han, Zhengyuan Zhou, Aaron Flores, and Tsachy Weissman. Leveraging the hints: Adaptive bidding in repeated first-price auctions. *NeurIPS*, 2022.

# A Appendix

## A.1 Measure and Information-Theoretic Notation and Known Facts

We recall that given two probability measures  $\mathbb{P}$  and  $\mathbb{Q}$  on a measurable space  $(\Omega, \mathcal{F})$ ,  $\mathbb{Q}$  is said to be absolutely continuous with respect to  $\mathbb{P}$  (and we write  $\mathbb{Q} \ll \mathbb{P}$ ) if, for all  $E \in \mathcal{F}$  such that  $\mathbb{P}[E] = 0$ , it holds that  $\mathbb{Q}[E] = 0$ . Whenever  $\mathbb{Q} \ll \mathbb{P}$ , the Radon-Nikodym theorem states that there exists a density (called Radon-Nikodym derivative of  $\mathbb{Q}$  with respect to  $\mathbb{P}$ )  $\frac{d\mathbb{Q}}{d\mathbb{P}}: \Omega \rightarrow [0, \infty)$  such that, for all  $E \in \mathcal{F}$ , it holds that

$$\mathbb{Q}[E] = \int_E \frac{d\mathbb{Q}}{d\mathbb{P}}(\omega) d\mathbb{P}(\omega).$$

See [Bass, 2013, Theorem 13.4] for a reference.

If  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space,  $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$  is a measurable space, and  $X$  is a random variable from  $(\Omega, \mathcal{F})$  to  $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ , the push-forward measure of  $\mathbb{P}$  by  $X$  is denoted by  $\mathbb{P}_X$ . In this case, we recall that the push-forward measure is defined as the unique probability measure on  $\mathcal{F}_{\mathcal{X}}$  defined via  $\mathbb{P}_X[F] = \mathbb{P}[X \in F]$ , for all  $F \in \mathcal{F}_{\mathcal{X}}$ .

If  $(\Omega, \mathcal{F})$  and  $(\Omega', \mathcal{F}')$  are two measurable spaces, their product  $\sigma$ -algebra is denoted by  $\mathcal{F} \otimes \mathcal{F}'$ . We recall that  $\mathcal{F} \otimes \mathcal{F}'$  is the  $\sigma$ -algebra of subsets of  $\Omega \times \Omega'$  generated by the collection of subsets of the form  $F \times F'$ , where  $F \in \mathcal{F}$  and  $F' \in \mathcal{F}'$ . If  $(\Omega, \mathcal{F}, \mathbb{P})$  and  $(\Omega', \mathcal{F}', \mathbb{P}')$  are two probability spaces, the product measure of  $\mathbb{P}$  and  $\mathbb{P}'$  is denoted by  $\mathbb{P} \otimes \mathbb{P}'$ . We recall that  $\mathbb{P} \otimes \mathbb{P}'$  is the unique probability measure defined on  $\mathcal{F} \otimes \mathcal{F}'$  which satisfies  $(\mathbb{P} \otimes \mathbb{P}')[F \times F'] = \mathbb{P}[F]\mathbb{P}'[F']$ , for all  $E \in \mathcal{F}$  and  $E' \in \mathcal{F}'$ . If  $(\Omega, \mathcal{F}, \mathbb{P})$  is a probability space,  $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$  and  $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$  are measurable spaces,  $X$  is a random variable from  $(\Omega, \mathcal{F})$  to  $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ , and  $Y$  is a random variable from  $(\Omega, \mathcal{F})$  to  $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ , the conditional probability of  $X$  given  $Y$  is denoted by  $\mathbb{P}_{X|Y}$ , where, for each  $E \in \mathcal{F}_{\mathcal{X}}$ , we recall that  $\mathbb{P}_{X|Y}[E] = \mathbb{P}[X \in E | Y]$  and that  $\mathbb{P}_{X|Y}[E]$  is a  $\sigma(Y)$ -measurable random variable.

The following result has been proven in Cesa-Bianchi et al. [2023].

**Theorem 9.** *Suppose that  $(\mathcal{Y}, d)$  is a separable and complete metric space with  $\mathcal{F}_{\mathcal{Y}}$  as the Borel  $\sigma$ -algebra of  $(\mathcal{Y}, d)$ . Let  $(\Omega, \mathcal{F})$  be a measurable space,  $X$  a random variable from  $(\Omega, \mathcal{F})$  to  $(\{0, 1\}, 2^{\{0, 1\}})$ ,  $Y$  a random variable from  $(\Omega, \mathcal{F})$  to  $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ , and  $U$  random variable from  $(\Omega, \mathcal{F})$  to  $([0, 1], \mathcal{B})$ , where  $\mathcal{B}$  is the Borel  $\sigma$ -algebra of  $[0, 1]$ . Suppose that  $\mathbb{P}, \mathbb{Q}$  are probability measures defined on  $\mathcal{F}$ , and  $p \in (0, 1)$ ,  $q \in [0, 1]$  are such that:*

- $\mathbb{P}[X = 1] = p$  and  $\mathbb{Q}[X = 1] = q$ .
- $U$  is a uniform random variable on  $[0, 1]$  both under  $\mathbb{P}$  and  $\mathbb{Q}$ , i.e., we have that  $\mathbb{P}_U = \mathbb{L} = \mathbb{Q}_U$ .
- $U$  is independent of  $X$  both under  $\mathbb{P}$  and  $\mathbb{Q}$ , i.e.,  $\mathbb{P}_{(X,U)} = \mathbb{P}_X \otimes \mathbb{P}_U$  and  $\mathbb{Q}_{(X,U)} = \mathbb{Q}_X \otimes \mathbb{Q}_U$ .

Then, the following are equivalent:

1. There exists a measurable function  $\varphi$  from  $(\{0, 1\} \times [0, 1], 2^{\{0, 1\}} \otimes \mathcal{B})$  to  $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$  such that

$$\mathbb{P}_Y = \mathbb{P}_{\varphi(X,U)} \quad \text{and} \quad \mathbb{Q}_Y = \mathbb{Q}_{\varphi(X,U)}.$$

2.  $\mathbb{Q}_Y \ll \mathbb{P}_Y$ , and  $\mathbb{P}_Y$ -almost-surely it holds that

$$\min \frac{d\mathbb{Q}_X}{d\mathbb{P}_X} \leq \frac{d\mathbb{Q}_Y}{d\mathbb{P}_Y} \leq \max \frac{d\mathbb{Q}_X}{d\mathbb{P}_X}.$$

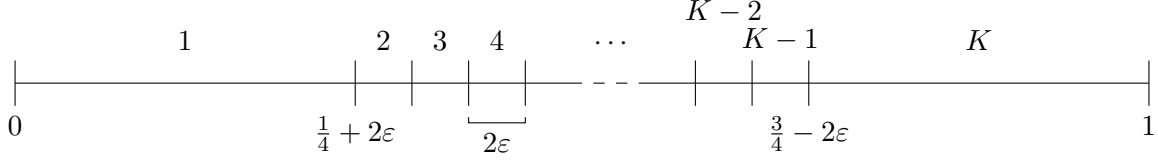


Figure 4: A representation of the map  $\iota$  through which the bids in the first-price auction problem are related to the  $K$ -arms of the bandit problem. The interval  $[0, 1]$  is partitioned in  $K$  disjoint intervals, the first and the last one of length  $\frac{1}{4} + 2\varepsilon$ , and all the ones in between of length  $2\varepsilon$ .  $\iota$  maps each bid to the index of the interval to which it belongs.

## A.2 Missing Details of the Proof of Theorem 3

In this section, we will complete the proof of Theorem 3, showing that the repeated first-price auctions with semi-transparent feedback (in the following, referred to as “our problem”) are no easier than a  $K$ -armed bandit instance based on the probability measures  $\mathbb{P}^1, \dots, \mathbb{P}^K$  introduced in Theorem 3. The structure of the proof is inspired by [Cesa-Bianchi et al., 2023, Section 3].

**The related bandit problem.** The action space is  $[K]$ , where we recall that  $K$  was some arbitrarily fixed natural number. Let  $Y, Y_1, Y_2, \dots$  be a sequence of  $\{0, 1\}^K$ -valued random variables such that, for any  $k \in \{0, 1, \dots, K\}$ , the sequence is  $\mathbb{P}^k$ -i.i.d. and, for all  $j \in [K]$

$$\mathbb{P}^k[Y(j) = 1] = \begin{cases} 1/2 & \text{if } j \neq k \\ 1/2 + 1/(6K) & \text{if } j = k \end{cases}$$

This sequence of latent random variables will determine the rewards of the actions. The reward function is

$$\rho: [K] \times \{0, 1\} \rightarrow [0, 1], \quad (i, y) \mapsto \frac{23 + 2y(i)}{192}$$

and the feedback received after playing an action  $I_t$  at time  $t$  is  $Y_t(I_t)$  (which is equivalent to receiving the bandit feedback  $\rho(I_t, Y_t)$  gathered at time  $t$ ).

For any  $k \in \{0, \dots, K\}$  and any  $i \in [K]$  the expected reward is

$$\mathbb{E}^k[\rho(i, Y)] = \begin{cases} \frac{1}{8} & \text{if } i \neq k \\ \frac{1}{8} + \frac{\varepsilon}{144} & \text{if } i = k \end{cases}$$

**Mapping our problem into this bandit problem.** Assume that  $K \geq 3$ . We partition the interval  $[0, 1]$  in the following  $K$  disjoint regions:  $J_1 = [0, w_1 + \varepsilon)$ ,  $J_k = [w_k - \varepsilon, w_k + \varepsilon)$  (for all  $k \in \{2, \dots, K - 1\}$ ), and  $J_K = [w_K - \varepsilon, 1]$ . We define a function  $\iota: [0, 1] \rightarrow [K]$  that maps each point in the interval  $[0, 1]$  to one of the  $K$  arms by mapping each  $b \in [0, 1]$  to the unique  $i \in [K]$  such that  $b \in J_i$  (for a pictorial representation of the map  $\iota$ , see Figure 4).

**Simulating the feedback.** To lighten the notation, besides the already defined random functions  $\psi_1, \psi_2, \dots$ , define also:

$$\psi: [0, 1] \rightarrow ([0, 1] \times \{\star\}) \cup (\{\star\} \times [0, 1]), \quad b \mapsto \begin{cases} (V, \star) & \text{if } b \geq M \\ (\star, M) & \text{if } b < M \end{cases}$$

The next lemma shows that we can use the feedback observed in the bandit problem together with some independent noise to simulate exactly the feedback of our problem.

**Lemma 6.** *For each  $b \in [0, 1]$ , there exists  $\varphi_b: \{0, 1\} \times [0, 1] \rightarrow ([0, 1] \times \{\star\}) \cup (\{\star\} \times [0, 1])$  such that, if  $U'$  is a  $[0, 1]$ -valued random variable such that, for each  $k \in \{0, \dots, K\}$ , the distribution  $U'$  with respect to  $\mathbb{P}^k$  is a uniform on  $[0, 1]$  and  $U'$  is  $\mathbb{P}^k$ -independent of  $Y$ , then  $\mathbb{P}_{\varphi_b(Y(\iota(b)), U')}^k = \mathbb{P}_{\psi(b)}^k$ .*

*Proof of Lemma 6.* A direct verification shows that, for all  $k \in [K]$  and all  $b \in [0, 1]$ ,  $\mathbb{P}_{\psi(b)}^k \ll \mathbb{P}_{\psi(b)}^0$  (i.e.,  $\mathbb{P}_{\psi(b)}^k$  is absolutely continuous with respect to  $\mathbb{P}_{\psi(b)}^0$ ) and the Radon-Nikodym derivative of the push-forward measure  $\mathbb{P}_{\psi(b)}^k$  with respect to  $\mathbb{P}_{\psi(b)}^0$  satisfies, for  $\mathbb{P}_{\psi(b)}^0$ -a.e.  $(v, m) \in ([0, 1] \times \{\star\}) \cup (\{\star\} \times [0, 1])$ ,

$$\frac{d\mathbb{P}_{\psi(b)}^k}{d\mathbb{P}_{\psi(b)}^0}(v, m) = 1 + \varepsilon \cdot \frac{16}{9} (v - b) \operatorname{sgn} \left( v - \frac{15}{16} \right) \Lambda_{w_k, \varepsilon}(b) \mathbb{I} \left\{ v \in \left[ \frac{7}{8}, 1 \right] \right\}$$

which implies, for  $\mathbb{P}_{\psi(b)}^0$ -a.e.  $(v, m) \in ([0, 1] \times \{\star\}) \cup (\{\star\} \times [0, 1])$ , that

$$\min \left( \frac{d\mathbb{P}_{Y(\iota(b))}^k}{d\mathbb{P}_{Y(\iota(b))}^0} \right) = 1 - \frac{4}{3} \varepsilon \leq \frac{d\mathbb{P}_{\psi(b)}^k}{d\mathbb{P}_{\psi(b)}^0}(v, m) \leq 1 + \frac{4}{3} \varepsilon = \max \left( \frac{d\mathbb{P}_{Y(\iota(b))}^k}{d\mathbb{P}_{Y(\iota(b))}^0} \right)$$

Thus, for each  $b \in [0, 1]$ , by Theorem 9, there exists (and we fix)

$$\varphi_b: \{0, 1\} \times [0, 1] \rightarrow ([0, 1] \times \{\star\}) \cup (\{\star\} \times [0, 1])$$

such that

$$\mathbb{P}_{\varphi_b(Y(\iota(b)), U')}^{\iota(b)} = \mathbb{P}_{\psi(b)}^{\iota(b)} \quad \text{and} \quad \mathbb{P}_{\varphi_b(Y(\iota(b)), U')}^0 = \mathbb{P}_{\psi(b)}^0.$$

Since for all  $b \in [0, 1]$  and all  $k \in [K] \setminus \{\iota(b)\}$ , we have  $\mathbb{P}_{\psi(b)}^k = \mathbb{P}_{\psi(b)}^0$  (by Equation (1)) and  $\mathbb{P}_{\varphi_b(Y(\iota(b)), U')}^k = \mathbb{P}_{\varphi_b(Y(\iota(b)), U')}^0$ , then, for all  $b \in [0, 1]$  and all  $k \in \{0, \dots, K\}$ , it holds that

$$\mathbb{P}_{\varphi_b(Y(\iota(b)), U')}^k = \mathbb{P}_{\psi(b)}^k. \quad \square$$

We now show that any algorithm  $\mathcal{A}$  for our problem can be transformed into an algorithm  $\tilde{\mathcal{A}}$  to solve the bandit problem that suffers no-larger regret. To do so, we begin by formally explaining how algorithms for our problem work.

**Functioning of an algorithm  $\mathcal{A}$  for our problem** A randomized algorithm  $\mathcal{A}$  for our problem is a sequence of functions that take as input a sequence of random seeds  $U_1, U_2, \dots$  and some feedback  $Z_1, Z_2, \dots$  and generates bids  $B_t$  as described below. At time  $t = 1$ ,  $\mathcal{A}$  selects a bid  $B_1$  as a deterministic function of  $U_1$  and observes feedback  $Z_1 = \psi_1(B_1)$ . Inductively, for any  $t \geq 2$ ,  $\mathcal{A}$  selects a bid  $B_t$  as a deterministic function of  $U_1, \dots, U_t, Z_1, \dots, Z_{t-1}$  (where  $Z_s = \psi_s(B_s)$ , for all  $s \in [t-1]$ ). For all  $k \in \{0, \dots, K\}$ , the sequence of seeds is a  $\mathbb{P}^k$ -i.i.d. sequence of uniform random variables on  $[0, 1]$  that is  $\mathbb{P}^k$ -independent of  $(V, M), (V_1, M_1), (V_2, M_2), \dots$ .

**Building  $\tilde{\mathcal{A}}$  from  $\mathcal{A}$**  We show now how to map  $\mathcal{A}$  to an algorithm  $\tilde{\mathcal{A}}$  (that shares the same seeds for the randomization) for the bandit problem that suffers a worst-case regret that is no larger than that of  $\mathcal{A}$ .

To do so, consider a sequence  $U', U'_1, \dots$  of random variables that, for all  $k \in \{0, \dots, K\}$  is a  $\mathbb{P}^k$ -i.i.d. sequence of uniforms on  $[0, 1]$  that  $\tilde{\mathcal{A}}$  can access as a further source of randomness. We will assume that, for all  $k \in \{0, \dots, K\}$ , the four sequences  $Y, Y_1, \dots, (V, M), (V_1, M_1), \dots, U, U_1, \dots$ , and  $U', U'_1, \dots$  are independent of each other.

The algorithm  $\tilde{\mathcal{A}}$  acts as follows. At time 1,  $\tilde{\mathcal{A}}$  plays the arm  $\tilde{I}_1 = \iota(B'_1)$ , where  $B'_1 = B_1$  is the bid played by  $\mathcal{A}$  at round  $t = 1$  (chosen as a deterministic function of the random seed  $U_1$ ). Then  $\tilde{\mathcal{A}}$  observes the bandit feedback  $Y_1(\tilde{I}_1)$  and feeds back to  $\mathcal{A}$  the surrogate feedback  $Z'_1 = \varphi_{B'_1}(Y_1(\tilde{I}_1), U'_1)$ . Then, inductively, for any time  $t \geq 2$ , assuming that  $\tilde{\mathcal{A}}$  played arms  $\tilde{I}_1, \dots, \tilde{I}_{t-1}$  and fed back to  $\mathcal{A}$  the surrogate feedback  $Z'_1, \dots, Z'_{t-1}$ , then

1.  $\tilde{\mathcal{A}}$  plays the arm  $\tilde{I}_t = \iota(B'_t)$ , where  $B'_t$  is the bid played by  $\mathcal{A}$  at round  $t$  (chosen as a deterministic function of the random seeds  $U_1, \dots, U_t$  and past surrogate feedback  $Z'_1, \dots, Z'_{t-1}$ ).
2.  $\tilde{\mathcal{A}}$  observes the bandit feedback  $Y_t(\tilde{I}_t)$  and feeds back to  $\mathcal{A}$  the surrogate feedback  $Z'_t = \varphi_{B'_t}(Y_t(\tilde{I}_t), U'_t)$ .

This way, we defined by induction the randomized algorithm  $\tilde{\mathcal{A}}$ .

By induction on  $t$ , one can show that, if  $B_1, B_2, \dots$  are the bids played by  $\mathcal{A}$  on the basis of the feedback  $Z_1 = \psi_1(B_1), Z_2 = \psi_2(B_2), \dots$ , then, for all  $k \in \{0, \dots, K\}$ , we have

$$\mathbb{P}_{(B_t, Y_t)}^k = \mathbb{P}_{(B'_t, Y_t)}^k$$

which leads to

$$\begin{aligned} R_T^k(\mathcal{A}) &= T \cdot \mathbb{E}^k[\text{Util}(w_k)] - \sum_{t=1}^T \mathbb{E}^k[\text{Util}_t(B_t)] \geq T \cdot \mathbb{E}^k[\rho(k, Y)] - \sum_{t=1}^T \mathbb{E}^k[\rho(\iota(B_t), Y_t)] \\ &= T \cdot \mathbb{E}^k[\rho(k, Y)] - \sum_{t=1}^T \mathbb{E}^k[\rho(\iota(B'_t), Y_t)] = T \cdot \mathbb{E}^k[\rho(k, Y)] - \sum_{t=1}^T \mathbb{E}^k[\rho(\tilde{I}_t, Y_t)] = \tilde{R}_T^k(\tilde{\mathcal{A}}) \end{aligned}$$

(the last equality is a definition). Now we are left to show only that for any algorithm  $\hat{\mathcal{A}}$  for the bandit problem which plays actions  $I_1, I_2, \dots$ , there exists  $k \in [K]$  such that

$$\tilde{R}_T^k(\hat{\mathcal{A}}) = T \cdot \mathbb{E}^k[\rho(k, Y)] - \sum_{t=1}^T \mathbb{E}^k[\rho(I_t, Y_t)] = \Omega(T^{2/3})$$

(the first equality is a definition). By Yao's Minimax principle, it is sufficient to show this for deterministic algorithms  $\hat{\mathcal{A}}$  for the bandit problem.

**Lemma 7.** *Fix any deterministic algorithm  $\hat{\mathcal{A}}$  for the bandit problem on  $K$  actions, then there exists  $k \in [K]$  such that  $\tilde{R}_T^k(\hat{\mathcal{A}}) \geq \frac{3}{10^4} T^{2/3}$ .*

*Proof.* For any deterministic algorithm  $\hat{\mathcal{A}}$  for the bandit problem on  $K$  actions, let  $I_1, I_2, \dots$  be the actions played by  $\hat{\mathcal{A}}$  on the basis of the sequential feedback received  $Z_1, Z_2, \dots$  and define  $N_t(i)$  as the random variables counting the number of times the learning algorithm  $\hat{\mathcal{A}}$  plays action  $i$ , up to time  $t$ , for any  $i \in [K]$  and any time  $t \in [T]$ :

$$N_t(i) = \sum_{s=1}^t \mathbb{I}\{I_s = i\}.$$

We relate the expected values of  $N_T(k)$  under  $\mathbb{P}^0$  and  $\mathbb{P}^k$  as a function of the expected number of times the algorithm plays the corresponding actions  $k$ . This formalizes the intuition that to discriminate between the different  $\mathbb{P}^k$  the learner needs to play exploring actions.

**Claim 3.** *The following inequality holds true for any  $k \in [K]$ :*

$$\mathbb{E}^k [N_T(k)] - \mathbb{E}^0 [N_T(k)] \leq \frac{2}{3} \cdot \varepsilon \cdot T \cdot \sqrt{2\mathbb{E}^0 [N_T(k)]}. \quad (5)$$

*Proof of Claim 3.* For any  $t \in [T]$ , the action  $I_t = I_t(Z_1, \dots, Z_{t-1})$  selected by  $\widehat{\mathcal{A}}$  at round  $t$  is a deterministic function of  $Z_1, \dots, Z_{t-1}$ , for each  $k \in [K]$ . In formula, we then have the following

$$\begin{aligned} \mathbb{E}^k [N_T(k)] - \mathbb{E}^0 [N_T(k)] &= \sum_{t=2}^T \left( \mathbb{P}^k [I_t(Z_1, \dots, Z_{t-1}) = k] - \mathbb{P}^0 [I_t(Z_1, \dots, Z_{t-1}) = k] \right) \\ &\leq \sum_{t=2}^T \left\| \mathbb{P}_{(Z_1, \dots, Z_{t-1})}^k - \mathbb{P}_{(Z_1, \dots, Z_{t-1})}^0 \right\|_{\text{TV}}, \end{aligned} \quad (6)$$

where  $\|\cdot\|_{\text{TV}}$  denotes the total variation norm. We move now our attention towards bounding the total variation norm. To that end we use Pinsker's inequality and apply the chain rule for the KL divergence KL. For each  $k \in [K]$  and  $t \in [T]$  we have the following:

$$\begin{aligned} \left\| \mathbb{P}_{(Z_1, \dots, Z_t)}^0 - \mathbb{P}_{(Z_1, \dots, Z_t)}^k \right\|_{\text{TV}} &\leq \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}_{(Z_1, \dots, Z_t)}^0, \mathbb{P}_{(Z_1, \dots, Z_t)}^k)} \\ &\leq \sqrt{\frac{1}{2} \left( \text{KL}(\mathbb{P}_{Z_1}^0, \mathbb{P}_{Z_1}^k) + \sum_{s=2}^t \mathbb{E} \left[ \text{KL}(\mathbb{P}_{Z_s | Z_1, \dots, Z_{s-1}}^0, \mathbb{P}_{Z_s | Z_1, \dots, Z_{s-1}}^k) \right] \right)} \end{aligned} \quad (7)$$

We bound the two KL terms separately.  $\widehat{\mathcal{A}}$  is a deterministic algorithm, thus  $I_1$  is a fixed element of  $[K]$ , which implies that, for all  $k \in [K]$ ,

$$\begin{aligned} &\text{KL}(\mathbb{P}_{Z_1}^0, \mathbb{P}_{Z_1}^k) \\ &= \left( \ln \left( \frac{\mathbb{P}^0 [Y_1(k) = 0]}{\mathbb{P}^k [Y_1(k) = 0]} \right) \mathbb{P}^0 [Y_1(k) = 0] + \ln \left( \frac{\mathbb{P}^0 [Y_1(k) = 1]}{\mathbb{P}^k [Y_1(k) = 1]} \right) \mathbb{P}^0 [Y_1(k) = 1] \right) \mathbb{I}\{I_1 = k\} \\ &= \frac{1}{2} \left( \ln \frac{1/2}{1/2 - \frac{2}{3} \cdot \varepsilon} + \ln \frac{1/2}{1/2 + \frac{2}{3} \cdot \varepsilon} \right) \cdot \mathbb{I}\{I_1 = k\} \end{aligned} \quad (8)$$

Similarly, since  $\widehat{\mathcal{A}}$  is a deterministic algorithm, for all  $s \geq 2$ , the action  $I_s = I_s(Z_1, \dots, Z_{s-1})$  selected

by  $\widehat{\mathcal{A}}$  at time  $t$  is a function of  $Z_1, \dots, Z_{s-1}$  only, which implies, for all  $k \in [K]$ ,

$$\begin{aligned}
& \text{KL}(\mathbb{P}_{Z_s|Z_1, \dots, Z_{s-1}}^0, \mathbb{P}_{Z_s|Z_1, \dots, Z_{s-1}}^k) \\
&= \mathbb{E}^0 \left[ \ln \left( \frac{\mathbb{P}^0[Z_s = 0 | Z_1, \dots, Z_{s-1}]}{\mathbb{P}^k[Z_s = 0 | Z_1, \dots, Z_{s-1}]} \right) \mathbb{P}^0[Z_s = 0 | Z_1, \dots, Z_{s-1}] \right. \\
&\quad \left. + \ln \left( \frac{\mathbb{P}^0[Z_s = 1 | Z_1, \dots, Z_{s-1}]}{\mathbb{P}^k[Z_s = 1 | Z_1, \dots, Z_{s-1}]} \right) \mathbb{P}^0[Z_s = 1 | Z_1, \dots, Z_{s-1}] \right] \\
&= \mathbb{E}^0 \left[ \left( \ln \left( \frac{\mathbb{P}^0[Y_s(k) = 0]}{\mathbb{P}^k[Y_s(k) = 0]} \right) \mathbb{P}^0[Y_s(k) = 0] + \ln \left( \frac{\mathbb{P}^0[Y_s(k) = 1]}{\mathbb{P}^k[Y_s(k) = 1]} \right) \mathbb{P}^0[Y_s(k) = 1] \right) \right. \\
&\quad \left. \times \mathbb{I}\{I_s(Z_1, \dots, Z_{s-1}) = k\} \right] \\
&= \frac{1}{2} \left( \ln \frac{1/2}{1/2 - \frac{2}{3} \cdot \varepsilon} + \ln \frac{1/2}{1/2 + \frac{2}{3} \cdot \varepsilon} \right) \mathbb{P}^0[I_s(Z_1, \dots, Z_{s-1}) = k] \tag{9}
\end{aligned}$$

Now, since  $\varepsilon = \frac{1}{4K} \leq \frac{1}{4} \leq \frac{2}{3}$ , the following useful inequality holds:

$$\frac{1}{2} \left( \ln \frac{1/2}{1/2 - \frac{2}{3} \cdot \varepsilon} + \ln \frac{1/2}{1/2 + \frac{2}{3} \cdot \varepsilon} \right) \leq 4 \cdot \left( \frac{2}{3} \right)^2 \cdot \varepsilon^2. \tag{10}$$

We can combine the inequalities in Equation (8) and Equation (9) into Equation (7) and plug in the bound in to obtain:

$$\left\| \mathbb{P}_{(Z_1, \dots, Z_t)}^0 - \mathbb{P}_{(Z_1, \dots, Z_t)}^k \right\|_{\text{TV}} \leq \frac{2}{3} \cdot \varepsilon \cdot \sqrt{2\mathbb{E}[N_t(k)]}$$

Once we have this upper bound on the total variations of the random variables  $(Z_1, \dots, Z_t)$  under  $\mathbb{P}^0$  and  $\mathbb{P}^k$  we can get back to the initial Equation (6) and obtain the desired bound via Jensen:

$$\mathbb{E}^k[N_T(k)] - \mathbb{E}^0[N_T(k)] \leq \sum_{t=2}^T \frac{2}{3} \cdot \varepsilon \cdot \sqrt{2\mathbb{E}^0[N_{t-1}(k)]} \leq \frac{2}{3} \cdot \varepsilon \cdot T \cdot \sqrt{2\mathbb{E}^0[N_T(k)]}. \quad \square$$

Averaging the quantitative bounds in Claim 3 for all  $k$  in  $[K]$ , and applying Jensen's inequality, we get the following:

$$\begin{aligned}
\frac{1}{K} \sum_{k \in [K]} \mathbb{E}^k[N_T(k)] &\leq \frac{1}{K} \sum_{k \in [K]} \mathbb{E}^0[N_T(k)] + \frac{2}{3} \cdot \varepsilon \cdot T \cdot \sqrt{\frac{2}{K} \sum_{k \in [K]} \mathbb{E}^0[N_T(k)]} \\
&= \left( \frac{1}{K} + \frac{2}{3} \cdot \varepsilon \cdot \sqrt{\frac{2T}{K}} \right) \cdot T. \tag{11}
\end{aligned}$$

Now, we have all the ingredients to lower bound the average regret suffered by  $\widehat{\mathcal{A}}$ . Note that every time a suboptimal arm is played the learner suffers (expected) instantaneous regret equal  $\frac{1}{144} \cdot \varepsilon$ .



Then, recalling that  $\varepsilon = 1/(4K)$  and setting  $K = \lceil T^{1/3} \rceil$  we have, for all  $T \geq 8$ ,

$$\begin{aligned} \frac{1}{K} \sum_{k \in [K]} \tilde{R}_T^k(\hat{\mathcal{A}}) &= \frac{1}{K} \sum_{k \in [K]} \left( \frac{1}{144} \cdot \varepsilon \cdot \mathbb{E}^k [T - N_T(k)] \right) = \frac{1}{144} \cdot \varepsilon \left( T - \frac{1}{K} \sum_{k \in [K]} \mathbb{E}^k [N_T(k)] \right) \\ &\geq \frac{1}{144} \cdot \varepsilon \cdot \left( 1 - \frac{1}{K} - \frac{2}{3} \cdot \varepsilon \cdot \sqrt{\frac{2T}{K}} \right) \cdot T = \frac{1}{144} \cdot \frac{1}{4K} \cdot \left( 1 - \frac{1}{K} - \frac{1}{6K} \cdot \sqrt{\frac{2T}{K}} \right) \cdot T \\ &\geq \frac{1}{8 \cdot 144} \left( \frac{3 - \sqrt{2}}{6} \right) T^{2/3} \geq \frac{3}{10^4} T^{2/3}. \end{aligned}$$

Therefore, for all  $T \geq 8$ , there exists  $k \in [K]$  such that  $\tilde{R}_T^k(\hat{\mathcal{A}}) \geq (3/10^4) \cdot T^{2/3}$ , concluding the proof.  $\square$

### A.3 Missing Details of the Proof of Theorem 5

**Claim 1.** *There exists two disjoint intervals  $I_+$  and  $I_-$  in  $[0, 1]$  such that, for any  $\varepsilon \in (0, 1/2)$  and any time  $t$ , the following hold:*

$$\max_{x \in [0, 1]} \mathbb{E}^{\pm\varepsilon} [\text{Util}_t(x)] \geq \mathbb{E}^{\pm\varepsilon} [\text{Util}_t(b)] + \frac{1}{128}\varepsilon, \text{ for all } b \notin I_{\pm}$$

*Proof.* For any  $\varepsilon \in (0, \frac{1}{2})$ , the distributions  $\mathbb{P}^{\pm\varepsilon}$  are such that, the set of all the bids that induce non-negative utility  $\mathbb{E}^{\pm\varepsilon} [\text{Util}_t(b)]$  is contained into two disjoint intervals  $I_+ = [0, \frac{1}{8}]$  and  $I_- = [\frac{1}{4}, 1]$ <sup>‡</sup>. We consider separately the two cases  $\mathbb{P}^{+\varepsilon}$  and  $\mathbb{P}^{-\varepsilon}$ . We start from the former. By simply looking at the definition (2), it is clear that  $\mathbb{E}^{+\varepsilon} [\text{Util}_t(b)]$  is monotonically increasing in  $\varepsilon$  for any  $b \in I_+$ , on the contrary, it is monotonically decreasing for  $b \in I_-$ . We have the following:

$$\max_{b \in I_-} \mathbb{E}^{+\varepsilon} [\text{Util}_t(b)] \leq \max_{b \in I_-} \mathbb{E}^0 [\text{Util}_t(b)] = \frac{1}{128}.$$

On the other hand,

$$\max_{x \in [0, 1]} \mathbb{E}^{+\varepsilon} [\text{Util}_t(x)] \geq \mathbb{E}^{+\varepsilon} [\text{Util}_t(\frac{1}{16})] = \frac{1}{128}(1 + \varepsilon) > \max_{b \in I_-} \mathbb{E}^{+\varepsilon} [\text{Util}_t(b)] + \frac{\varepsilon}{128}.$$

We consider now the other case, corresponding to  $\mathbb{P}^{-\varepsilon}$ . By the definition in Equation (2),  $\mathbb{E}^{-\varepsilon} [\text{Util}_t(b)]$  is monotonically increasing in its first argument for any  $b \in I_-$ , on the contrary, it is monotonically decreasing for  $b \in I_+$ . Similarly to the other case we have two steps. On the one hand, it holds that

$$\max_{b \in I_+} \mathbb{E}^{-\varepsilon} [\text{Util}_t(b)] \leq \max_{b \in I_+} \mathbb{E}^0 [\text{Util}_t(b)] = \frac{1}{128},$$

while on the other hand it holds that

$$\max_{x \in [0, 1]} \mathbb{E}^{-\varepsilon} [\text{Util}_t(x)] \geq \mathbb{E}^{-\varepsilon} [\text{Util}_t(\frac{7}{16})] = \frac{1}{128} + \varepsilon \frac{41}{128} > \max_{b \in I_-} \mathbb{E}^{-\varepsilon} [\text{Util}_t(b)] + \frac{\varepsilon}{4}. \quad \square$$

We need a preliminary result for the proof of Claim 2. Recall, we use the same random variable  $(V, M)$  to denote the highest competing bid/valuation pair drawn from the different probability distribution. When we change the underlying measure, we are changing its law. Consider now the push forward measures on  $[0, 1]^2$  (with the Borel  $\sigma$ -algebra) induced by these three measures:  $\mathbb{P}_{(V, M)}^0$ ,  $\mathbb{P}_{(V, M)}^{+\varepsilon}$  and  $\mathbb{P}_{(V, M)}^{-\varepsilon}$ . With some simple calculations (similarly to what is done in, e.g., Appendix B of Slivkins [2019]) it is possible to bound the KL divergence:

<sup>‡</sup>The choice of  $I_+$  and  $I_-$  is not tight.

**Claim 4.** For any  $\varepsilon \in (0, \frac{1}{2})$  the following inequality holds true:

$$\text{KL} \left( \mathbb{P}_{(V,M)}^{+\varepsilon}, \mathbb{P}_{(V,M)}^0 \right) = \text{KL} \left( \mathbb{P}_{(V,M)}^{-\varepsilon}, \mathbb{P}_{(V,M)}^0 \right) \leq 2\varepsilon^2$$

*Proof.* We simply apply the definition of KL divergence for continuous random variables. We only do the calculations for  $\mathbb{P}_{(V,M)}^{+\varepsilon}$ , the other term is analogous:

$$\begin{aligned} \text{KL} \left( \mathbb{P}_{(V,M)}^{+\varepsilon}, \mathbb{P}_{(V,M)}^0 \right) &= \int_{Q_+ \cup Q_-} f^{+\varepsilon}(v, m) \ln \frac{f^{+\varepsilon}(v, m)}{f^0(v, m)} dm dv \\ &= \frac{1}{2}(1 + \varepsilon) \ln(1 + \varepsilon) + \frac{1}{2}(1 - \varepsilon) \ln(1 - \varepsilon) \leq 2\varepsilon^2, \end{aligned}$$

where the last inequality holds for any  $\varepsilon \in (0, \frac{1}{2})$ . □

**Claim 2.** The following inequality holds:  $\frac{1}{2} \sum_{i=1,2} \mathbb{E}^i [N_i] \leq \frac{3}{4}T$ .

*Proof.* We have the following:

$$\begin{aligned} \mathbb{E}^i [N_i] - \mathbb{E}^0 [N_i] &= \sum_{t=2}^T \mathbb{P}^i [B_t \in I_i] - \mathbb{P}^0 [B_t \in I_i] \\ &\leq \sum_{t=2}^T \left\| \mathbb{P}_{(V_1, M_1), \dots, (V_{t-1}, M_{t-1})}^i - \mathbb{P}_{(V_1, M_1), \dots, (V_{t-1}, M_{t-1})}^0 \right\|_{\text{TV}} \quad (\text{Total variation}) \\ &\leq \sum_{t=2}^T \sqrt{\frac{1}{2} \text{KL} \left( \mathbb{P}_{(V_1, M_1), \dots, (V_{t-1}, M_{t-1})}^i, \mathbb{P}_{(V_1, M_1), \dots, (V_{t-1}, M_{t-1})}^0 \right)} \quad (\text{Pinsker's inequality}) \\ &\leq \sum_{t=2}^T \sqrt{\frac{t}{2} \text{KL} \left( \mathbb{P}_{(V,M)}^i, \mathbb{P}_{(V,M)}^0 \right)} \quad ((V_1, M_1), \dots, (V_{t-1}, M_{t-1}), \dots \text{ are i.i.d.}) \\ &\leq \frac{1}{4\sqrt{T}} \sum_{t=2}^T \sqrt{t} \leq \frac{1}{4}T, \end{aligned} \tag{12}$$

where in the last inequality we applied Claim 4 for our choice of  $\varepsilon = 1/(4\sqrt{T})$ . Note,  $\mathbb{P}_{(V_1, M_1), \dots, (M_t, V_t)}^j$  is the push-forward measure on  $([0, 1]^2)^t$  induced by  $t$  i.i.d. draws of  $(V, M)$  from distribution  $\mathbb{P}^j$ ,  $j \in \{0, 1, 2\}$ . Averaging the result in Equation (12), we get the desired inequality:

$$\frac{1}{2} \sum_{i=1,2} \mathbb{E}^i [N_i] \leq \frac{1}{2} \sum_{i=1,2} \mathbb{E}^0 [N_i] + \frac{T}{4} = \frac{3}{4}T. \quad \square$$