

AN ONLINE LEARNING THEORY OF TRADING-VOLUME MAXIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

We explore brokerage between traders in an online learning framework. At any round t , two traders meet to exchange an asset, provided the exchange is mutually beneficial. The broker proposes a trading price, and each trader tries to sell their asset or buy the asset from the other party, depending on whether the price is higher or lower than their private valuations. A trade happens if one trader is willing to sell and the other is willing to buy at the proposed price. Previous work provided guidance to a broker aiming at enhancing traders' total earnings by maximizing the *gain from trade*, defined as the sum of the traders' net utilities after each interaction. This classical notion of reward can be highly unfair to traders with small profit margins, and far from the real-life utility of the broker. For these reasons, we investigate how the broker should behave to maximize the trading volume, i.e., the *total number of trades*. We model the traders' valuations as an i.i.d. process with an unknown distribution. If the traders' valuations are revealed after each interaction (full-feedback), and the traders' valuations cumulative distribution function (cdf) is continuous, we provide an algorithm achieving logarithmic regret and show its optimality up to constants. If only their willingness to sell or buy at the proposed price is revealed after each interaction (2-bit feedback), we provide an algorithm achieving poly-logarithmic regret when the traders' valuations cdf is Lipschitz and show its near-optimality. We complement our results by analyzing the implications of dropping the regularity assumptions on the unknown traders' valuations cdf. If we drop the continuous cdf assumption, the regret rate degrades to $\Theta(\sqrt{T})$ in the full-feedback case, where T is the time horizon. If we drop the Lipschitz cdf assumption, learning becomes impossible in the 2-bit feedback case.

1 INTRODUCTION

In modern financial markets, Over-the-Counter (OTC) trading platforms have emerged as dynamic and decentralized hubs, offering diverse alternatives to traditional exchanges. In recent years, these markets have experienced remarkable growth, solidifying their central role in the global financial ecosystem: OTC asset trading in the US surpassed 50 trillion USD in value in 2020 (Weill, 2020), with an upward trend documented since 2016 (www.bis.org, 2022).

Brokers play a crucial role in OTC markets. Beyond acting as intermediaries between traders, they utilize their understanding of the market to identify the optimal prices for assets. Additionally, traders in these markets often respond to price changes: higher prices usually lead to selling, while lower prices typically result in buying (Sherstyuk et al., 2020). This adaptability appears across various asset classes, including stocks, derivatives, art, collectibles, precious metals and minerals, energy commodities (like gas and oil), and digital currencies (cryptocurrencies) (Bolić et al., 2024).

Our study draws inspiration from recent research analyzing the bilateral trade problem from an online learning perspective (Cesa-Bianchi et al., 2021; Azar et al., 2022; Cesa-Bianchi et al., 2023; 2024a; Bolić et al., 2024; Bernasconi et al., 2024; Bachoc et al., 2024a;b). In particular, we build on insights from Bolić et al. (2024), which addresses the brokerage problem in OTC markets where traders may decide to buy or sell their assets depending on prevailing market conditions.

054 1.1 MOTIVATIONS FOR CHOOSING TRADING VOLUME AS REWARD

055
056 Previous works have entirely focused on scenarios where brokers aim at maximizing the so-called
057 cumulative *gain from trade*—the sum of the net utilities of the traders over the entire sequence of
058 interactions with the broker. This classical approach has the two following pitfalls.

059
060 **Traders’ Perspective.** Gain-from-trade maximization can cause unfairness in settings where the
061 majority of traders make a living off of small margins (e.g., in micro trading or high-frequency
062 trading), and only a handful of high-payoff trades have the potential to occur. In these cases, gain-
063 from-trade maximization can lead to sacrificing the majority of the population in favor of a small
064 minority of traders that are lucky enough to be paired with people that are willing to be greatly
065 underpaid for the good on sale. In contrast, trading-volume maximization gives the same dignity to
066 all traders, granting everybody the same opportunity to trade, independently of their buying power.
067 For a striking concrete example of this pitfall, see Section 3.

068 **Broker’s Perspective.** From the broker’s perspective, too, it might not be as beneficial to potentially
069 miss out on traders’ exchanges by maximizing the gain from trade, given that, typically, brokers only
070 earn when trades occur. For example, in settings where traders have to pay a small fee for each trade,
071 it is clear that the broker’s ultimate goal is to maximize trading volume. Another example where
072 maximizing trading volume is superior to maximizing the gain from trade is the one discussed in the
073 Trader’s Perspective paragraph (and Section 3). In this case, a gain-from-trade maximizing broker
074 would risk alienating the vast majority of the population which, realistically, would end up leaving a
075 broker that does not give them trading opportunities, consequently hurting the broker’s bottom line.

076 For these reasons, in this work, we aim at providing strategies that boost the trading volume by
077 maximizing the *number of trades* in the broker-traders interaction sequence.

078 1.2 SETTING

079
080 In what follows, for any two real numbers a, b , we denote their minimum by $a \wedge b$ and their maximum
081 by $a \vee b$. We now describe the brokerage online learning protocol.

082
083 For any time $t = 1, 2, \dots$

- 084 • Two traders arrive with their private valuations V_{2t-1} and V_{2t}
- 085 • The broker proposes a trading price P_t
- 086 • If the price P_t is between the lowest valuation $V_{2t-1} \wedge V_{2t}$ and the highest valuation $V_{2t-1} \vee$
087 V_{2t} —meaning the trader with the lower valuation is willing to sell at P_t and the trader with
088 the higher valuation is willing to buy at P_t —the transaction occurs with the higher-valuation
089 trader purchasing the asset from the lower-valuation trader at the price P_t
- 090 • The broker receives some feedback

091
092 As commonly assumed in the existing bilateral trade literature, we assume valuations and prices
093 belong to $[0, 1]$. While previous literature aims at maximizing the cumulative *gain from trade*—
094 defined as the sum of traders’ net utilities¹ in the whole interaction sequence—our objective is to
095 maximize the *number of trades*. Formally, for any $p, v_1, v_2 \in [0, 1]$, our utility posting a price p when
096 the valuations of the traders are v_1 and v_2 is

$$097 \quad g(p, v_1, v_2) := \mathbb{I}\{v_1 \wedge v_2 \leq p \leq v_1 \vee v_2\} .$$

098 The goal of the broker is to minimize the *regret*, defined, for any time horizon $T \in \mathbb{N}$, as

$$099 \quad R_T := \sup_{p \in [0, 1]} \mathbb{E} \left[\sum_{t=1}^T (G_t(p) - G_t(P_t)) \right] ,$$

100
101 where $G_t(q) := g(q, V_{2t-1}, V_{2t})$ for all $q \in [0, 1]$ and $t \in \mathbb{N}$, and the expectation is taken over the
102 randomness present in $(V_t)_{t \in \mathbb{N}}$ and the (possible) randomness used by the broker’s algorithm to
103 generate the prices $(P_t)_{t \in \mathbb{N}}$.

104
105 ¹Formally, for any $p, v_1, v_2 \in [0, 1]$, the gain from trade of a price p when the valuations of the traders are v_1
106 and v_2 is $\text{GFT}(p, v_1, v_2) := (v_1 \vee v_2 - v_1 \wedge v_2) \mathbb{I}\{v_1 \wedge v_2 \leq p \leq v_1 \vee v_2\}$.

	M -Lipschitz	Continuous	General
Full	$\Omega(\ln T)$ Thm 2	$O(\ln T)$ Thm 1	$\Theta(\sqrt{T})$ Thm 5+6
2-Bit	$O(\ln(MT) \ln T), \Omega(\ln(MT))$ Thm 3+4	$\Omega(T)$ Thm 7	$\Omega(T)$

Table 1: Overview of all the regret regimes: $\ln T$ (cyan), $\ln(MT)$ (green), \sqrt{T} (yellow), and T (red), depending on the feedback (full or 2-bit) and the assumption on the cdf (M -Lipschitz, continuous, or no assumptions).

As in Bolić et al. (2024), we assume that traders’ valuations V, V_1, V_2, \dots are generated i.i.d. from an *unknown* distribution ν —a practical assumption for large and stable markets.

Finally, we consider the following two different types of feedback commonly studied in the online learning bilateral trade literature:

- *Full-feedback*. At each round t , after having posted the price P_t , the broker has access to the traders’ valuations V_{2t-1} and V_{2t} .
- *2-bit feedback*. At each round t , after having posted the price P_t , the broker has access to the indicator functions $\mathbb{I}\{V_{2t-1} \leq P_t\}$ and $\mathbb{I}\{V_{2t} \leq P_t\}$.

The full-feedback model draws its motivation from *direct revelation mechanisms*, where the traders disclose their valuations V_{2t-1} and V_{2t} before each round, but the mechanism has access to this information only after having posted the current bid P_t (Cesa-Bianchi et al., 2021; 2024a).

The 2-bit feedback model corresponds to *posted price mechanisms*, where the broker has access only to the traders’ willingness to buy or sell at the proposed posted price, and the valuations V_{2t-1} and V_{2t} are *never* revealed.

1.3 OVERVIEW OF OUR CONTRIBUTIONS

In the full-feedback case, if the distribution ν of the traders’ valuations has a *continuous* cdf, we design an algorithm (Algorithm 1) suffering $O(\ln T)$ regret in the time horizon T (Theorem 1), and we provide a matching lower bound (Theorem 2). We complement these results by showing that dropping the continuous cdf assumption leads to a worse regret rate of $\Omega(\sqrt{T})$ (Theorem 5), and we design an algorithm (Algorithm 3) achieving $O(\sqrt{T})$ regret (Theorem 6).

In the 2-bit feedback case, if the cdf of the traders’ valuations is M -Lipschitz, we design an algorithm (Algorithm 2) achieving regret $O(\ln(MT) \ln T)$ (Theorem 3) where T is the time horizon, and provide a near-matching lower bound $\Omega(\ln(MT))$ (Theorem 4). We complement these results by showing that the problem becomes unlearnable if we drop the Lipschitzness assumption (Theorem 7).

For a full summary of our results, see Table 1.

1.4 TECHNIQUES AND CHALLENGES

Online learning with a continuous action domain and full-feedback is usually tackled by discretizing the action domain and then playing an optimal expert algorithm on the discretization, or by directly running exponential weights algorithms in the continuum (Maillard & Munos, 2010; Krichene et al., 2015; Cesa-Bianchi et al., 2024b). These approaches require that the (expected) reward function is Lipschitz and lead to a regret rate of order $\tilde{O}(\sqrt{T})$. In contrast, our expected reward function is *not* Lipschitz in general. To overcome this challenge, we leverage the specific structure of the problem by proving Lemma 1, which enables us to design an algorithm that achieves an exponentially better regret rate of $O(\ln T)$ even when the underlying cdf—and hence the associated reward function—is only continuous. Moreover, we establish a matching $\Omega(\ln T)$ lower bound that, surprisingly, applies even when the reward function is Lipschitz, demonstrating that additional Lipschitz regularity beyond continuity does not contribute to faster rates in this setting. This lower bound construction is particularly challenging because the shape of the function $p \mapsto \mathbb{E}[G_t(p)]$ can only be controlled indirectly through the traders’ valuation distribution: to avoid exceedingly complex calculations, extra care is required in selecting appropriate instances. Even then, we needed a subtle and somewhat intricate Bayesian argument to obtain the lower bound.

In the 2-bit feedback model, we remark that the available feedback is enough to reconstruct *bandit* feedback. Consequently, when the underlying cdf—and hence the expected reward function—is M -Lipschitz, a viable approach is to discretize the action space $[0, 1]$ with K uniformly spaced points and run an optimal bandit algorithm on the discretization. This approach immediately yields a regret rate of order $O(MT/K + \sqrt{KT})$. This bound leads to a regret of order $O(M^{1/3}T^{2/3})$ by tuning $K := \Theta(M^{2/3}T^{1/3})$ when M is known to the learner, or of order $O(MT^{2/3})$ by tuning $K := \Theta(T^{2/3})$ when the learner does not possess this knowledge. In contrast, we exploit the extra information provided by the 2-bit feedback and the intuition provided by Lemma 1 to devise a binary search algorithm achieving the exponentially better rate of $O(\ln(MT) \ln T)$, with the additional feature of being oblivious to M . Our corresponding lower bound shows that this rate is optimal (up to a $\ln T$ factor), demonstrating through an information-theoretic argument that some sort of binary search is essentially a necessary step for optimal learning.

1.5 RELATED WORK

Since the pioneering work of Myerson and Satterthwaite (Myerson & Satterthwaite, 1983), the study of bilateral trade has grown significantly, particularly from a game-theoretic and approximation perspective (Colini-Baldeschi et al., 2016; 2017; Blumrosen & Mizrahi, 2016; Brustle et al., 2017; Colini-Baldeschi et al., 2020; Babaioff et al., 2020; Dütting et al., 2021; Deng et al., 2022; Kang et al., 2022; Archbold et al., 2023). For a comprehensive overview, refer to Cesa-Bianchi et al. (2024a).

In recent years, the focus has expanded to include online learning settings for bilateral trade. Given their close relevance to our paper, we concentrate our discussion on these works.

In Cesa-Bianchi et al. (2021); Azar et al. (2022); Cesa-Bianchi et al. (2024a; 2023); Bernasconi et al. (2024); Cesa-Bianchi et al. (2024b), the authors examined bilateral trade problems where the reward function is the *gain from trade* and each trader has a fixed role as either a seller or a buyer.

In Cesa-Bianchi et al. (2021), the authors investigated a scenario where seller and buyer valuations form two distinct i.i.d. sequences. In the full-feedback case, they achieved a regret bound of $\tilde{O}(\sqrt{T})$, which was later refined to $O(\sqrt{T})$ in Cesa-Bianchi et al. (2024a). They also demonstrated a worst-case regret of $\Omega(\sqrt{T})$ even when sellers’ and buyers’ valuations are independent of each other and their cdfs are Lipschitz. For the 2-bit feedback scenario under i.i.d. valuations, Cesa-Bianchi et al. (2021) proved that any algorithm must suffer linear regret, even under either the M -Lipschitz joint cdf assumption or the traders’ valuation independence assumption. However, when both conditions are simultaneously satisfied, they proposed an algorithm achieving a regret rate of $\tilde{O}(M^{1/3}T^{2/3})$, later refined to $O(M^{1/3}T^{2/3})$ in Cesa-Bianchi et al. (2024a). Cesa-Bianchi et al. (2021) also established a worst-case regret lower bound of $\Omega(T^{2/3})$ in this case, which, however, does not display any dependence on M .

Cesa-Bianchi et al. (2021; 2024a) also showed that the adversarial bilateral trade problem is unlearnable even with full-feedback. To achieve learnability beyond the i.i.d. case, Cesa-Bianchi et al. (2023; 2024b) explored weakly budget-balanced mechanisms, allowing the broker to post different selling and buying prices as long as the buyer pays more than what the seller receives. They demonstrated that learning can be achieved using weakly budget-balanced mechanisms in the 2-bit feedback setting at a regret rate of $\tilde{O}(MT^{3/4})$ when the joint seller/buyer cdf may vary over time but is M -Lipschitz. Furthermore, for the same setting, they provided a $\Omega(T^{3/4})$ matching lower bound in the time horizon, even when the process is required to be i.i.d., but their lower bound does not feature any dependence on M . Azar et al. (2022) investigated whether learning is possible in the adversarial case by considering α -regret, achieving $\tilde{\Theta}(\sqrt{T})$ bounds for 2-regret in full-feedback, and a $\tilde{O}(T^{3/4})$ upper bound in 2-bit feedback. Following another direction, Bernasconi et al. (2024) explored globally budget-balanced mechanisms in the adversarial case, showing a $\Theta(\sqrt{T})$ regret rate in full-feedback and a $\tilde{O}(T^{3/4})$ rate in the 2-bit feedback case.

The closest to our work is Bolić et al. (2024), where the authors studied the same i.i.d. version of our trading problem with flexible seller and buyer roles, but with the target reward function being the *gain from trade*. Under the M -Lipschitz cdf assumption, they obtained tight $\Theta(M \ln T)$ regret in the full-feedback case. Surprisingly, in the same full-feedback case, but using our different reward function, we achieve a $\Theta(\ln T)$ regret rate even when the cdf is only continuous: in our case, the

216 additional Lipschitz regularity does not offer any speedup once the continuity assumption is fulfilled.
 217 Furthermore, under the M -Lipschitz cdf assumption, Bolić et al. (2024) proved a sharp rate of
 218 $\Theta(\sqrt{MT})$ in the 2-bit feedback case. Interestingly, using our different reward function, we achieve
 219 an exponentially faster upper bound of $O(\ln(MT) \ln T)$, which is tight up to a $\ln T$ factor. If the
 220 Lipschitz cdf assumption is removed, the learning rate for both our problem and the one in Bolić et al.
 221 (2024) degrades to $\Theta(\sqrt{T})$ in the full-feedback case, and the problem becomes unlearnable in the
 222 2-bit feedback case.

224 2 THE MEDIAN LEMMA

226 In this section, we present the Median Lemma (Lemma 1), a simple but crucial result for what follows,
 227 and the key upon which our main algorithms are based. At a high level, Lemma 1 states that a broker
 228 who aims at maximizing the number of trades should post prices that are as close as possible to the
 229 *median* of the (unknown) traders' valuation distribution ν , and the instantaneous regret which the
 230 broker incurs by playing any price is (proportional to) the *square* of the distance between the median
 231 and the price, if distances are measured with respect to the pseudo-metric induced by the traders'
 232 valuation cdf.

233 **Lemma 1** (The median lemma). *If the cdf F of ν is continuous, then, for any $t \in \mathbb{N}$ and any $p \in [0, 1]$,*

$$235 \mathbb{E}[G_t(p)] = 2F(p)(1 - F(p)) \quad \text{and} \quad \frac{1}{2} - \mathbb{E}[G_t(p)] = 2\left(\frac{1}{2} - F(p)\right)^2.$$

237 *In particular, the function $p \mapsto \mathbb{E}[G_t(p)]$ is maximized at any point $m \in [0, 1]$ such that $F(m) = \frac{1}{2}$.*

239 Before presenting the proof of Lemma 1, we just remark that points $m \in [0, 1]$ satisfying $F(m) = 1/2$
 240 do exist by the intermediate value theorem, because $F(0) = 0$, $F(1) = 1$, and F is continuous.

242 *Proof.* For each $t \in \mathbb{N}$ and each $p \in [0, 1]$, we have that

$$243 \mathbb{E}[G_t(p)] = \mathbb{P}\left[\{V_{2t-1} \leq p < V_{2t}\} \cup \{V_{2t} \leq p \leq V_{2t-1}\}\right]
 244 = \mathbb{P}[V_{2t-1} \leq p] \mathbb{P}[p < V_{2t}] + \mathbb{P}[V_{2t} \leq p] \mathbb{P}[p \leq V_{2t-1}] = 2F(p)(1 - F(p)),$$

247 where the second equality follows from additivity and independence, while in the last equality we
 248 leveraged the continuity of F to obtain $\mathbb{P}[p \leq V_{2t-1}] = \mathbb{P}[p < V_{2t-1}] = 1 - F(p)$. To conclude, it is
 249 enough to note that, for each $p \in [0, 1]$ it holds that $\frac{1}{4} - F(p)(1 - F(p)) = (\frac{1}{2} - F(p))^2$. \square

251 3 TRADING VOLUME VS GAIN FROM TRADE

253 In this section, we leverage Lemma 1 to show with a formal example that, unlike trading-volume
 254 maximizing brokers, gain-from-trade maximization brokers can be heavily biased towards small
 255 segments of the population and, as a result, end up hurting their own bottom lines.

256 Assume that the distribution of the traders' valuations V, V_1, V_2, \dots have common density f defined,
 257 for all $x \in [0, 1]$, by $f(x) := (\frac{1}{\varepsilon} - 1) \mathbb{I}\{\frac{1}{2} - \varepsilon \leq x \leq \frac{1}{2}\} + \mathbb{I}\{1 - \varepsilon \leq x \leq 1\}$, for some $\varepsilon \in (0, \frac{1}{2})$.

259 At a high level, this population of traders is clustered into two segments: a *low*-valuation cluster L
 260 that believes that the good on sale has a value slightly smaller than $\frac{1}{2}$ and a *high*-valuation cluster H
 261 that believes the value is slightly smaller than 1. If $\varepsilon \approx 0$, the overwhelming majority of the population
 262 belongs to the low-valuation cluster L . In this case, we will prove that a gain-from-trade maximizing
 263 broker would sacrifice the majority of the population to favor trades that include a trader coming
 264 from the (extremely small) high-valuation cluster H .

265 Indeed, by Bolić et al. (2024), a gain-from-trade maximizing broker would post the *expectation*
 266 $\mathbb{E}[V] = \frac{1}{2}$. In contrast, by Lemma 1, a trade-volume maximizing broker would post the *median*
 267 $m := \frac{1}{2} - \frac{\varepsilon}{2} \cdot \frac{1-2\varepsilon}{1-\varepsilon}$ of V , which is a value roughly in the middle of the low-valuation cluster L .

268 By posting the expectation, the probability of having a trade is, for all $t \in \mathbb{N}$, $\mathbb{P}[V_{2t-1} \wedge V_{2t} \leq \frac{1}{2} \leq$
 269 $V_{2t-1} \vee V_{2t}] = 2(1 - \varepsilon)\varepsilon$, which is close to zero when $\varepsilon \approx 0$.

In contrast, by posting the median, the probability of having a trade is, for all $t \in \mathbb{N}$, $\mathbb{P}[V_{2t-1} \wedge V_{2t} \leq m \leq V_{2t-1} \vee V_{2t}] = \frac{1}{2}$, which is always bounded away from zero, irrespectively of ε .

This shows two ways in which (unlike a trade-volume maximizing broker) a gain-from-trade maximizing broker is biased towards favoring the high valuation cluster H . First, they are willing to accept that only a negligible fraction of the population will trade. Second, being $\mathbb{E}[V] = \frac{1}{2}$, they only (with probability 1) allow trades where one of the two traders comes from the high-valuation cluster H , resulting in *only* the high-valuation trader making a large profit, while the low valuation trader is left with a profit of order $\varepsilon \approx 0$, even in the low-probability event where they are given the opportunity to trade. It is easy to imagine that, in real life, such a bias would cause the low-valuation traders in L to leave the broker, in turn greatly reducing the broker's own profit.

4 FULL-FEEDBACK

We now investigate how the broker should behave to maximize the number of trades in the full-feedback case where after each interaction the traders' valuations are disclosed. We begin by studying the full-feedback case under the continuous cdf assumption. In this case, taking inspiration from Lemma 1, a natural strategy is to play the *empirical median*, which leads to Algorithm 1.

Algorithm 1: Follow the Empirical Median (FEM)

Post $P_1 := 1/2$ and receive feedback V_1, V_2 ;

for time $t = 2, 3, \dots$ **do**

Post the empirical median $P_t := \frac{1}{2} (V_{2(t-1)}^{(t-1)} + V_{2(t-1)}^{(t)})$, where $V_{2(t-1)}^{(1)}, \dots, V_{2(t-1)}^{(2(t-1))}$ are the order statistics of the observed sample $V_1, \dots, V_{2(t-1)}$, and receive feedback V_{2t-1}, V_{2t} ;

The next theorem leverages Lemma 1 to show that Algorithm 1 suffers regret $O(\ln T)$ when the traders' valuation cdf is continuous.

Theorem 1. *If ν has a continuous cdf F , the regret of FEM satisfies, for all time horizons $T \in \mathbb{N}$,*

$$R_T \leq \frac{1}{2} + \frac{\pi}{2} (1 + \ln(T-1)).$$

Proof. Without loss of generality, we can (and do!) assume that $T \geq 2$. Then, we have

$$R_T \leq \frac{1}{2} + \max_{p \in [0,1]} \mathbb{E} \left[\sum_{t=2}^T G_t(p) \right] - \mathbb{E} \left[\sum_{t=2}^T G_t(P_t) \right] = \frac{1}{2} + 2 \cdot \sum_{t=2}^T \mathbb{E} \left[\left(\frac{1}{2} - F(P_t) \right)^2 \right]$$

Now, let $m \in [0, 1]$ be such that $F(m) = 1/2$, and let V be a random variable whose distribution is ν , independent of V_1, V_2, \dots . Then, for any $t \in \mathbb{N}$ such that $t \geq 2$ we have

$$\mathbb{E} \left[\left(\frac{1}{2} - F(P_t) \right)^2 \right] = \mathbb{E} \left[(\mathbb{P}[m \leq V \leq P_t | P_t])^2 \right] + \mathbb{E} \left[(\mathbb{P}[P_t \leq V \leq m | P_t])^2 \right] =: (I) + (II).$$

Now, for the term (I) , leveraging the fact that V and P_t are independent of each other, together with the Minkowski's integral inequality (see, e.g., (Stein, 1970, Appendix A.1)), we have:

$$\begin{aligned} \sqrt{(I)} &= \sqrt{\mathbb{E} \left[(\mathbb{E}[\mathbb{I}\{m \leq V \leq P_t\} | P_t])^2 \right]} \leq \mathbb{E} \left[\sqrt{\mathbb{E} \left[(\mathbb{I}\{m \leq V \leq P_t\})^2 | V \right]} \right] \\ &= \mathbb{E} \left[\sqrt{\mathbb{P}[m \leq V \leq P_t | V]} \right] = \int_{[m,1]} \sqrt{\mathbb{P}[x \leq P_t]} d\mathbb{P}_V(x) = \int_{[m,1]} \sqrt{\mathbb{P}[x \leq P_t]} d\nu(x) = (\star) \end{aligned}$$

For each $x \in [0, 1]$ and for any $s \in \mathbb{N}$, let $B_s(x) := \mathbb{I}\{x \leq V_s\}$, and notice that $B_1(x), B_2(x), \dots$ is an i.i.d. sequence of Bernoulli random variables of parameter $1 - F(x)$. Let $V_{2(t-1)}^{(1)}, \dots, V_{2(t-1)}^{(2(t-1))}$ be the order statistics of the observed sample $V_1, \dots, V_{2(t-1)}$. For any $x \in [m, 1]$, observing that

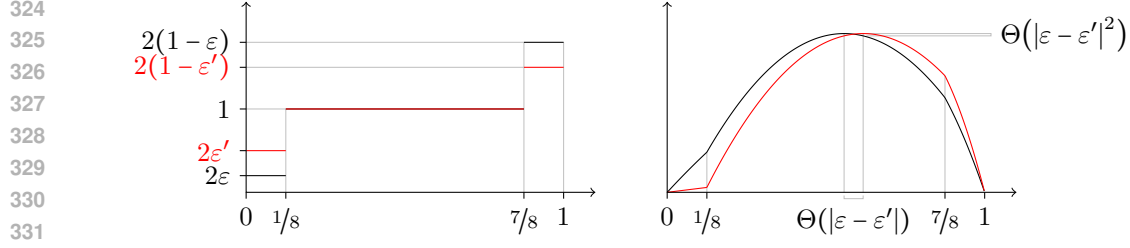


Figure 1: Qualitative plots of the densities $f_\varepsilon, f_{\varepsilon'}$ (left) and corresponding expected rewards (right) used in the proof of Theorem 2 for two values $\varepsilon, \varepsilon' > 0$. Note that the difference in reward by posting a price that is optimal for one instance ε' when the actual instance is ε is $\Theta(|\varepsilon - \varepsilon'|^2)$.

$F(x) - \frac{1}{2} \geq 0$ and $\mathbb{P}[x \leq P_t] \leq \mathbb{P}[x \leq V_{2(t-1)}^{(t)}] \leq \mathbb{P}\left[\sum_{s=1}^{2(t-1)} B_s(x) \geq t-1\right]$, we can leverage Hoeffding's inequality to obtain

$$\begin{aligned} \mathbb{P}[x \leq P_t] &\leq \mathbb{P}\left[\sum_{s=1}^{2(t-1)} B_s(x) \geq t-1\right] = \mathbb{P}\left[\sum_{s=1}^{2(t-1)} \frac{B_s(x)}{2(t-1)} - (1-F(x)) \geq \frac{t-1}{2(t-1)} - (1-F(x))\right] \\ &= \mathbb{P}\left[\sum_{s=1}^{2(t-1)} \frac{B_s(x)}{2(t-1)} - (1-F(x)) \geq F(x) - \frac{1}{2}\right] \leq e^{-4(t-1)(F(x)-\frac{1}{2})^2} = e^{-4(t-1)(\nu[[0,x]]-\frac{1}{2})^2}, \end{aligned}$$

from which, by the change of variable formula (Revuz & Yor, 2013, Proposition 4.10, Chapter 1), it follows also that

$$\begin{aligned} (*) &\leq \int_{[m,1]} \sqrt{\exp\left(-4(t-1)\left(\nu[[0,x]] - \frac{1}{2}\right)^2\right)} d\nu(x) = \int_{1/2}^1 \exp\left(-2(t-1)\left(\frac{1}{2}-u\right)^2\right) du \\ &\leq \frac{1}{\sqrt{2(t-1)}} \int_0^\infty \exp(-r^2) dr = \frac{\sqrt{\pi}}{2\sqrt{2}} \cdot \frac{1}{\sqrt{t-1}}, \end{aligned}$$

and hence $(I) \leq \frac{\pi}{8(t-1)}$. Analogously, we can prove that $(II) \leq \frac{\pi}{8(t-1)}$. Hence,

$$R_T \leq \frac{1}{2} + \frac{\pi}{2} \cdot \sum_{t=2}^T \frac{1}{t-1} = \frac{1}{2} + \frac{\pi}{2} + \frac{\pi}{2} \cdot \sum_{t=2}^{T-1} \int_{t-1}^t \frac{1}{t} ds \leq \frac{1}{2} + \frac{\pi}{2} + \frac{\pi}{2} \cdot \int_1^{T-1} \frac{1}{s} ds = \frac{1}{2} + \frac{\pi}{2} (1 + \ln(T-1)). \quad \square$$

We now establish the optimality of FEM by demonstrating a matching $\Omega(\ln T)$ regret lower bound. We remark that this result holds even when competing against underlying distributions that have a 2-Lipschitz cdf, thus proving the optimality of FEM even under the Lipschitz cdf assumption.

Theorem 2. *There exist two numerical constants c_1 and c_2 such that, for any time horizon $T \geq c_2$, the worst-case regret of any full-feedback algorithm satisfies*

$$\sup_{\nu \in \mathcal{D}_2} R_T^\nu \geq c_1 \ln T,$$

where R_T^ν is the regret at time T of the algorithm when the i.i.d. sequence of traders' valuations follows the distribution ν , and \mathcal{D}_2 is the set of all distributions ν that admit a 2-Lipschitz cdf.

Due to space constraints, we defer the (long and technical) proof of this result to Appendix A and only present a short, high-level sketch here.

Proof sketch. In the proof, we build a family of 2-Lipschitz cdfs F_ε parameterized by $\varepsilon \in [0, 1]$, so that if two instances are parameterized by ε and ε' respectively, then their medians are $\Theta(|\varepsilon - \varepsilon'|)$ -away from each other (Figure 1). The high-level idea is to leverage a Bayesian argument to show that if the underlying instance F_E is such that E is drawn uniformly at random in $[0, 1]$, then, at round t , the broker cannot reliably determine prices that are much closer than $1/\sqrt{t}$ to the corresponding median m_E when distances are measured with respect to the metric induced by the cdf F_E . This, together with our key Lemma 1, leads to the conclusion. \square

5 2-BIT FEEDBACK

We start the study of the 2-bit feedback case under the assumption that the traders' valuation distribution admits a Lipschitz cdf F . The algorithm we propose (Algorithm 2) is based on the following observation: by posting any price p , the broker has access to two noisy realizations of $F(p)$. Recalling that Lemma 1 suggests tracking the median of F (i.e., a point m at which $F(m) = 1/2$), and since F is a non-decreasing function, we can proceed using a natural binary search strategy to move toward the median. This can be done in epochs: in each one, we repeatedly test a (dyadic) price until the first time we can confidently decide that the median is to the left or right of the current price.

Algorithm 2: Median Binary Search (MBS)

Input: Confidence parameter $\delta \in (0, 1)$, time horizon $n \in \mathbb{N}$;

Initialization: $Q_1 := \frac{1}{2}$, $\tau := 1$, $t := 1$;

while time $t \leq n$ **do**

Let $s := 0$, $Y_{\tau,s} := 0$, $t_{\tau-1} := t - 1$;

while time $t \leq n$ **do**

Post $P_t := Q_\tau$ and receive feedback $\mathbb{I}\{V_{2t-1} \leq P_t\}$, $\mathbb{I}\{V_{2t} \leq P_t\}$;

Update $s := s + 2$, $Y_{\tau,s} := Y_{\tau,s-2} + \mathbb{I}\{V_{2t-1} \leq P_t\} + \mathbb{I}\{V_{2t} \leq P_t\}$, $t := t + 1$;

if $\frac{1}{s}Y_{\tau,s} + \sqrt{\frac{\ln(2/\delta)}{2s}} < \frac{1}{2}$ **then** let $Q_{\tau+1} := Q_\tau + \frac{1}{2^{\tau+1}}$, $s_\tau := s$, $\tau := \tau + 1$, and **break**;

else if $\frac{1}{s}Y_{\tau,s} - \sqrt{\frac{\ln(2/\delta)}{2s}} > \frac{1}{2}$ **then** let $Q_{\tau+1} := Q_\tau - \frac{1}{2^{\tau+1}}$, $s_\tau := s$, $\tau := \tau + 1$, and **break**;

We now show that a suitably tuned Algorithm 2 has regret guarantees of $O(\ln(MT) \ln(T))$. In particular, we stress that the tuning of Algorithm 2 does not need prior knowledge of M . Due to space constraints, we defer the full proof of the next result to Appendix B.

Theorem 3. *If ν has an M -Lipschitz cdf F for some $M \geq 1$, then, for all time horizons $T \in \mathbb{N}$, the regret of MBS tuned with parameters $\delta := 2/T^3$ and $n := T$ satisfies*

$$R_T \leq 2 + 6 \log_2(MT) \ln(T).$$

Proof sketch. The proof is based on the following observations. First, during an epoch where a price p is tested, given that one has to distinguish if the parameter $F(p)$ of a sequence of Bernoulli random variables is bigger or smaller than $1/2$, a concentration argument shows that the duration of this epoch is at most $O(\ln(1/\delta)/(1/2 - F(p))^2)$, where δ is the confidence parameter. Second, by Lemma 1, the broker regrets $2(1/2 - F(p))^2$ by playing a price p , and hence the total regret of an epoch where the broker tests p is at most $O(\ln(1/\delta))$. We then use the fact that the F is M -Lipschitz to prove that, after at most $O(\log_2(MT))$ epochs, the cumulative regret that the algorithm suffers from that point onward is constant, and conclude by showing that the tuning of δ leads to the stated guarantees. \square

We now show that Algorithm 2 is optimal, up to a multiplicative $\ln T$ term. Due to space constraints, we defer the full proof of this result to Appendix C.

Theorem 4. *There exist two numerical constants c_1 and c_2 such that for any $M \geq 16$ and any time horizon $T \geq c_2 \log_2(M)$, the worst-case regret of any 2-bit feedback algorithm satisfies*

$$\sup_{\nu \in \mathcal{D}_M} R_T^\nu \geq c_1 \ln(MT),$$

where R_T^ν is the regret at time T of the algorithm when the i.i.d. sequence of traders' valuations follows the distribution ν , and \mathcal{D}_M is the set of all distributions ν that admits an M -Lipschitz cdf.

Proof sketch. The proof builds a family of distributions, each supported in a different region of length $\Theta(1/M)$, whose cdfs are M -Lipschitz. To avoid suffering linear regret if the traders' valuations are generated according to one of these distributions, the broker has to detect the corresponding support. To accomplish this task, we show that the broker is essentially forced to solve a binary search problem that needs $\log_2(M)$ rounds in each of which the instantaneous regret is constant. Noticing that any

regret lower bound for full-feedback algorithms also applies to 2-bit feedback algorithms, the $\ln T$ lower bound of Theorem 2 together with the binary search $\ln M$ lower bound yield a lower bound of $\Omega(\max(\ln T, \ln M)) = \Omega(\ln(MT))$. \square

6 NON-LIPSCHITZ OR DISCONTINUOUS PDFS

We now investigate how the problem changes if we lift the assumption that ν has a Lipschitz or continuous cdf. First, note that when the cdf of ν is not continuous, Lemma 1, and, consequently, the guarantees of Theorem 1, no longer hold. Indeed, in general, no full-feedback algorithm can achieve regret guarantees better than \sqrt{T} . As shown in the proof of the next theorem, the reason is that our problem contains instances that resemble online learning with expert advice (with 2 experts), which has a known lower bound of $\Omega(\sqrt{T})$.

Theorem 5. *There exist two numerical constants c_1 and c_2 such that, for any time horizon $T \geq c_2$, the worst-case regret of any full feedback algorithm satisfies*

$$\sup_{\nu \in \mathcal{D}} R_T^\nu \geq c_1 \sqrt{T},$$

where R_T^ν is the regret at time T of the algorithm when the i.i.d. sequence of traders' valuations follows the distribution ν , and \mathcal{D} is the set of all distributions ν .

Proof sketch. For each $\varepsilon \in [-1/4, 1/4]$, define $\nu_\varepsilon := \frac{1-\varepsilon}{4}\delta_0 + \frac{1}{4}\delta_{1/3} + \frac{1}{4}\delta_{2/3} + \frac{1+\varepsilon}{4}\delta_1$, where, for any $a \in \mathbb{R}$, we denoted by δ_a the Dirac's delta probability measure centered at a . Let $(V_{\varepsilon,t})_{\varepsilon \in [-1/4, 1/4], t \in \mathbb{N}}$ be an independent family such that for each $\varepsilon \in [-1/4, 1/4]$ the sequence $V_{\varepsilon,1}, V_{\varepsilon,2}, \dots$ is i.i.d. with common distribution ν_ε . For each $\varepsilon \in [-1/4, 1/4]$, each $t \in \mathbb{N}$, and each $p \in [0, 1]$, define $G_{\varepsilon,t}(p) := g(p, V_{\varepsilon,2t-1}, V_{\varepsilon,2t})$. Straightforward computations show that, for each $\varepsilon \in [-1/4, 1/4]$ and each $t \in \mathbb{N}$, the function $p \mapsto \mathbb{E}[G_{\varepsilon,t}(p)]$ is maximized at $1/3$ or at $2/3$, with any other point having an expected reward that is less than $31/256$ -away from the minimum expected reward achieved at $1/3$ or $2/3$. Furthermore, for any $\varepsilon \in [-1/4, 1/4]$ and any $t \in \mathbb{N}$, the maximum is at $1/3$ or $2/3$ depending on whether $\varepsilon < 0$ or $\varepsilon > 0$, given that $\mathbb{E}[G_{\varepsilon,t}(1/3)] = \frac{11}{16} - \frac{\varepsilon}{8} - \frac{\varepsilon^2}{8}$ and $\mathbb{E}[G_{\varepsilon,t}(2/3)] = \frac{11}{16} + \frac{\varepsilon}{8} - \frac{\varepsilon^2}{8}$, from which it follows also that $\mathbb{E}[G_{\varepsilon,t}(2/3)] - \mathbb{E}[G_{\varepsilon,t}(1/3)] = \frac{\varepsilon}{4}$. Hence, in order not to suffer $\Omega(|\varepsilon|T)$ regret, an algorithm has to detect the *sign* of ε . However, a standard information-theoretic argument shows that a sample of order $\Omega(1/\varepsilon^2)$ is required in order to detect the sign of ε . During this period, the best any algorithm can do is to play blindly in the set $\{1/3, 2/3\}$, incurring in a cumulative regret of order $\Omega(\frac{1}{\varepsilon^2} \cdot |\varepsilon|) = \Omega(1/|\varepsilon|)$. Overall, any learner has to suffer $\Omega(\min(\frac{1}{|\varepsilon|}, |\varepsilon|T))$ worst-case regret, which, by tuning $|\varepsilon| = \Theta(1/\sqrt{T})$, leads to a worst-case regret lower bound of $\Omega(\sqrt{T})$. \square

We now focus on the upper bound. A closer look at the proof of Lemma 1 shows that if we drop the cdf continuity assumption in the Median Lemma, the formula generalizes to

$$\mathbb{E}[G_t(p)] = 2F(p)(1 - F(p)) + F(p)F^\circ(p) =: \Psi(p), \quad \forall p \in [0, 1], \forall t \in \mathbb{N},$$

with no assumptions on ν , and where F is the cdf of ν and we defined $F^\circ(p) := \nu[\{p\}]$. This suggests the strategy of building an empirical proxy $\hat{\Psi}_t$ of Ψ with the feedback available at time t , and posting prices that maximize $\hat{\Psi}_t$. By replacing the theoretical quantities by their empirical counterparts, for any $t \in \mathbb{N}$ and any $p \in [0, 1]$, we define an empirical proxy for $\Psi(p)$ as follows:

$$\hat{\Psi}_{t+1}(p) := 2 \frac{1}{2t} \sum_{s=1}^{2t} \mathbb{I}\{V_s \leq p\} \frac{1}{2t} \sum_{s=1}^{2t} \mathbb{I}\{p < V_s\} + \frac{1}{2t} \sum_{s=1}^{2t} \mathbb{I}\{V_s \leq p\} \frac{1}{2t} \sum_{s=1}^{2t} \mathbb{I}\{V_s = p\}.$$

This definition leads to Algorithm 3.

We now state regret guarantees for Algorithm 3. The proof of the following result (which hinges on showing that $\hat{\Psi}_t$ is *uniformly* close to Ψ with high probability) is deferred to Appendix D.

Theorem 6. *For all time horizons $T \in \mathbb{N}$, the regret of FE Ψ satisfies*

$$R_T \leq 1 + 8\sqrt{\pi} \cdot \sqrt{T-1}.$$

Algorithm 3: Follow the Empirical Ψ (FE Ψ)

Post $P_1 = \frac{1}{2}$ and receive feedback V_1, V_2 ;

for time $t = 2, 3, \dots$ **do**

Post $P_t \in \operatorname{argmax}_{p \in [0,1]} \hat{\Psi}_t(p)$ and receive feedback V_{2t-1}, V_{2t} ;

We conclude by showing that, without the Lipschitz cdf assumption, the 2-bit feedback problem is unlearnable. This can be deduced as a simple corollary of the proof of Theorem 4. Specifically, we can obtain a linear worst-case lower bound for any 2-bit feedback algorithm, even if we assume that the underlying distribution has a continuous cdf.

Theorem 7. *There exist two numerical constants c_1 and c_2 such that, for any time horizon $T \geq c_2$, the worst-case regret of any 2-bit feedback algorithm satisfies*

$$\sup_{\nu \in \mathcal{D}_c} R_T^\nu \geq c_1 T,$$

where R_T^ν is the regret at time T of the algorithm when the i.i.d. sequence of traders' valuations follows the distribution ν , and \mathcal{D}_c is the set of all distributions ν that admits a continuous cdf.

Proof. As a consequence of the last part of the proof of Theorem 4 (see Appendix C) we have that, for any time horizon $T \geq 4$, if we set $M := 2^T$, then the conditions $M \geq 16$ and $T \geq \log_2(M)$ in that proof holds, and hence, any 2-bit feedback algorithm has worst-case regret that is at least $\frac{1}{4 \ln 2} \ln M = \frac{1}{4 \ln 2} T$. \square

7 CONCLUSIONS AND OPEN PROBLEMS

Motivated by maximizing trading volume in OTC markets, we proposed a novel objective that departs from the classical *gain-from-trade* reward studied in the bilateral trade literature. For this new problem, we investigated optimal brokerage strategies from an online learning perspective. Under the assumption that traders are free to sell or buy depending on the trading price and that traders' valuations form an i.i.d. sequence, we provided a complete picture with matching (up to, at most, logarithmic factors) upper and lower bounds in all the proposed settings, fleshing out the role of regularity assumptions in achieving these fast regret rates.

In addition to closing the logarithmic $\ln T$ gap in the regret rate of the 2-bit feedback setting, a few other future research directions are to find non-stationary variants of this problem where learning is still achievable, investigate trading volume maximization when traders have definite seller and buyer roles, and explore the contextual version of the problem when the broker has access to relevant side information before posting each price.

REFERENCES

- Thomas Archbold, Bart de Keijzer, and Carmine Ventre. Non-obvious manipulability for single-parameter agents and bilateral trade. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 2107–2115, USA, 2023. International Foundation for Autonomous Agents and Multiagent Systems.
- Yossi Azar, Amos Fiat, and Federico Fusco. An alpha-regret analysis of adversarial bilateral trade. *Advances in Neural Information Processing Systems*, 35:1685–1697, 2022.
- Moshe Babaioff, Kira Goldner, and Yannai A. Gonczarowski. Bulow-klemperer-style results for welfare maximization in two-sided markets. In *Proceedings of the Thirty-First Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '20, pp. 2452–2471, USA, 2020. Society for Industrial and Applied Mathematics.
- François Bachoc, Nicolò Cesa-Bianchi, Tommaso Cesari, and Roberto Colomboni. Fair online bilateral trade. *arXiv preprint arXiv:2405.13919*, 2024a.

- 540 François Bachoc, Tommaso Cesari, and Roberto Colomboni. A contextual online learning theory of
541 brokerage. *arXiv preprint arXiv:2407.01566*, 2024b.
- 542
- 543 Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. No-regret learning in
544 bilateral trade via global budget balance. In *Proceedings of the 56th Annual ACM Symposium on*
545 *Theory of Computing*, 2024.
- 546 Liad Blumrosen and Yehonatan Mizrahi. Approximating gains-from-trade in bilateral trading. In
547 *Web and Internet Economics, WINE'16*, volume 10123 of *Lecture Notes in Computer Science*, pp.
548 400–413, Germany, 2016. Springer.
- 549 Natasa Bolić, Tommaso Cesari, and Roberto Colomboni. An online learning theory of brokerage. In
550 *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*,
551 AAMAS '24, pp. 216–224, Richland, SC, 2024. International Foundation for Autonomous Agents
552 and Multiagent Systems. ISBN 9798400704864.
- 553
- 554 Johannes Brustle, Yang Cai, Fa Wu, and Mingfei Zhao. Approximating gains from trade in two-sided
555 markets via simple mechanisms. In *Proceedings of the 2017 ACM Conference on Economics*
556 *and Computation, EC '17*, pp. 589–590, New York, NY, USA, 2017. Association for Computing
557 Machinery. ISBN 9781450345279.
- 558 Nicolò Cesa-Bianchi, Tommaso R Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi.
559 A regret analysis of bilateral trade. In *Proceedings of the 22nd ACM Conference on Economics*
560 *and Computation*, pp. 289–309, USA, 2021. Association for Computing Machinery.
- 561 Nicolò Cesa-Bianchi, Tommaso R Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi.
562 Repeated bilateral trade against a smoothed adversary. In *The Thirty Sixth Annual Conference on*
563 *Learning Theory*, pp. 1095–1130, USA, 2023. PMLR, PMLR.
- 564
- 565 Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi.
566 Bilateral trade: A regret minimization perspective. *Mathematics of Operations Research*, 49(1):
567 171–203, 2024a.
- 568 Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi.
569 Regret analysis of bilateral trade with a smoothed adversary. *Journal of Machine Learning*
570 *Research*, 25(234):1–36, 2024b.
- 571
- 572 Tommaso R Cesari and Roberto Colomboni. A nearest neighbor characterization of Lebesgue points
573 in metric measure spaces. *Mathematical Statistics and Learning*, 3(1):71–112, 2021.
- 574 Riccardo Colini-Baldeschi, Bart de Keijzer, Stefano Leonardi, and Stefano Turchetta. Approximately
575 efficient double auctions with strong budget balance. In *ACM-SIAM Symposium on Discrete*
576 *Algorithms, SODA'16*, pp. 1424–1443, USA, 2016. SIAM.
- 577
- 578 Riccardo Colini-Baldeschi, Paul W. Goldberg, Bart de Keijzer, Stefano Leonardi, and Stefano
579 Turchetta. Fixed price approximability of the optimal gain from trade. In *Web and Internet*
580 *Economics, WINE'17*, volume 10660 of *Lecture Notes in Computer Science*, pp. 146–160, Germany,
581 2017. Springer.
- 582 Riccardo Colini-Baldeschi, Paul W Goldberg, Bart de Keijzer, Stefano Leonardi, Tim Roughgar-
583 den, and Stefano Turchetta. Approximately efficient two-sided combinatorial auctions. *ACM*
584 *Transactions on Economics and Computation (TEAC)*, 8(1):1–29, 2020.
- 585 Yuan Deng, Jieming Mao, Balasubramanian Sivan, and Kangning Wang. Approximately efficient
586 bilateral trade. In *STOC*, pp. 718–721, Italy, 2022. ACM.
- 587
- 588 Paul Dütting, Federico Fusco, Philip Lazos, Stefano Leonardi, and Rebecca Reiffenhäuser. Efficient
589 two-sided markets with limited information. In *Proceedings of the 53rd Annual ACM SIGACT*
590 *Symposium on Theory of Computing*, STOC 2021, pp. 1452–1465, New York, NY, USA, 2021.
591 Association for Computing Machinery. ISBN 9781450380539.
- 592 Zi Yang Kang, Francisco Pernice, and Jan Vondrák. Fixed-price approximations in bilateral trade.
593 In *Proceedings of the 2022 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp.
2964–2985, Alexandria, VA, USA, 2022. SIAM, Society for Industrial and Applied Mathematics.

- 594 Walid Krichene, Maximilian Balandat, Claire Tomlin, and Alexandre Bayen. The Hedge algorithm
595 on a continuum. In *International Conference on Machine Learning*, pp. 824–832. PMLR, 2015.
596
- 597 Odalric-Ambrym Maillard and Rémi Munos. Online learning in adversarial lipschitz environments.
598 In *ECML/PKDD (2)*, volume 6322 of *Lecture Notes in Computer Science*, pp. 305–320. Springer,
599 2010.
- 600 Pascal Massart. The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The annals of*
601 *Probability*, pp. 1269–1283, 1990.
602
- 603 Roger B Myerson and Mark A Satterthwaite. Efficient mechanisms for bilateral trading. *Journal of*
604 *economic theory*, 29(2):265–281, 1983.
- 605 Daniel Revuz and Marc Yor. *Continuous martingales and Brownian motion*, volume 293. Springer
606 Science & Business Media, 2013.
607
- 608 Katerina Sherstyuk, Krit Phankitnirundorn, and Michael J Roberts. Randomized double auctions:
609 gains from trade, trader roles, and price discovery. *Experimental Economics*, 24(4):1–40, 2020.
610
- 611 Elias M Stein. *Singular integrals and differentiability properties of functions*. Princeton university
612 press, 1970.
- 613 Pierre-Olivier Weill. The search theory of over-the-counter markets. *Annual Review of Economics*,
614 12:747–773, 2020.
- 615 www.bis.org. OTC derivatives statistics at end-June 2022. *Bank for International Settlements*, 2022.
616 URL https://www.bis.org/publ/otc_hy2211.pdf.
617

619 A PROOF OF THEOREM 2

620 For each $\varepsilon \in [0, 1]$, consider the following density function (see Figure 1, left)

$$621 f_\varepsilon: [0, 1] \rightarrow [0, 2], \quad x \mapsto 2\varepsilon \mathbb{I}\left\{x \leq \frac{1}{8}\right\} + \mathbb{I}\left\{\frac{1}{8} < x < \frac{7}{8}\right\} + 2(1 - \varepsilon) \mathbb{I}\left\{x \geq \frac{7}{8}\right\},$$

622 Notice that, for each $\varepsilon \in [0, 1]$ the cumulative function associated to the density f_ε is 2-Lipschitz
623 with explicit expression given by

$$624 F_\varepsilon: [0, 1] \rightarrow [0, 1], \quad x \mapsto 2\varepsilon x \mathbb{I}\left\{x \leq \frac{1}{8}\right\} + \left(\frac{2\varepsilon - 1}{8} + x\right) \mathbb{I}\left\{\frac{1}{8} < x < \frac{7}{8}\right\} + (2\varepsilon - 1 - 2(\varepsilon - 1)x) \mathbb{I}\left\{x \geq \frac{7}{8}\right\}.$$

625 Consider for each $\varepsilon \in [0, 1]$, an i.i.d. sequence $(B_{\varepsilon,t})_{t \in \mathbb{N}}$ of Bernoulli random variables of parameter
626 ε , an i.i.d. sequence $(D_t)_{t \in \mathbb{N}}$ of Bernoulli random variables of parameter $\frac{1}{4}$, an i.i.d. sequence
627 $(U_t)_{t \in \mathbb{N}}$ of uniform random variables on $[0, 1]$, and a uniform random variable E on $[0, 1]$, such
628 that $((B_{\varepsilon,t})_{t \in \mathbb{N}, \varepsilon \in [0, 1]}, (D_t)_{t \in \mathbb{N}}, (U_t)_{t \in \mathbb{N}}, E)$ is an independent family. For each $\varepsilon \in [0, 1]$ and $t \in \mathbb{N}$,
629 define

$$630 V_{\varepsilon,t} := D_t \cdot \left(B_{\varepsilon,t} \frac{U_t}{8} + (1 - B_{\varepsilon,t}) \frac{7 + U_t}{8} \right) + (1 - D_t) \cdot \left(\frac{1}{8} + \frac{3}{4} U_t \right). \quad (1)$$

631 Tedious but straightforward computations show that, for each $\varepsilon \in [0, 1]$ the sequence $(V_{\varepsilon,t})_{t \in \mathbb{N}}$ is
632 i.i.d. with common density given by f_ε , and this sequence is independent of E . For any $\varepsilon \in [0, 1]$,
633 $p \in [0, 1]$, and $t \in \mathbb{N}$, let $G_{\varepsilon,t}(p) := g(p, V_{\varepsilon,2t-1}, V_{\varepsilon,2t})$ (for a qualitative representation of its
634 expectation, see Figure 1, right). We now show how to lower bound the worst-case regret of any
635 arbitrary deterministic algorithm for the full-feedback setting $(\alpha_t)_{t \in \mathbb{N}}$, i.e., a sequence of functions
636 $\alpha_t: ([0, 1] \times [0, 1])^{t-1} \rightarrow [0, 1]$ where each element maps past feedback into a price (with the
637 convention that α_1 is a number in $[0, 1]$). We remark that we do not lose any generality in considering
638 only deterministic algorithms given that we are in a stochastic i.i.d. setting, and the minimax regret
639 over deterministic algorithms coincides with that over randomized algorithms. For each $t \in \mathbb{N}$, define
640 $\tilde{\alpha}_t: ([0, 1] \times [0, 1])^{t-1} \rightarrow [\frac{1}{8}, \frac{7}{8}]$ equal to α_t whenever α_t takes values in $[\frac{1}{8}, \frac{7}{8}]$, and equal to $1/2$
641 otherwise. Notice that for each $\varepsilon \in [0, 1]$ it holds that $(F_\varepsilon \circ \tilde{\alpha}_t) \cdot (1 - F_\varepsilon \circ \tilde{\alpha}_t) \geq (F_\varepsilon \circ \alpha_t) \cdot (1 - F_\varepsilon \circ \alpha_t)$,
642 and hence, due to Lemma 1, for each $t \in \mathbb{N}$, it holds that $\mathbb{E}[G_{\varepsilon,t}(\tilde{\alpha}_t(V_{\varepsilon,1}, \dots, V_{\varepsilon,2(t-1)}))] \geq$

648 $\mathbb{E} \left[G_{\varepsilon, t}(\alpha_t(V_{\varepsilon, 1}, \dots, V_{\varepsilon, 2(t-1)})) \right]$. Notice also that for each $\varepsilon \in [0, 1]$, we have that $m_\varepsilon := \frac{5-2\varepsilon}{8}$ is the
 649 unique element in $[0, 1]$ such that $F_\varepsilon(m_\varepsilon) = 1/2$. For any time horizon $T \geq 144$, we have that the
 650 worst-case regret of the algorithm $(\alpha_t)_{t \in \mathbb{N}}$ can be lower bounded as follows

$$\begin{aligned}
 651 & \\
 652 & \sup_{\nu \in \mathcal{D}_M} R_T^\nu \geq \sup_{\varepsilon \in [0, 1]} \sum_{t=13}^T \mathbb{E} \left[G_{\varepsilon, t}(m_\varepsilon) - G_{\varepsilon, t}(\alpha_t(V_{\varepsilon, 1}, \dots, V_{\varepsilon, 2(t-1)})) \right] \\
 653 & \\
 654 & \geq \sup_{\varepsilon \in [0, 1]} \sum_{t=13}^T \mathbb{E} \left[G_{\varepsilon, t}(m_\varepsilon) - G_{\varepsilon, t}(\tilde{\alpha}_t(V_{\varepsilon, 1}, \dots, V_{\varepsilon, 2(t-1)})) \right] \stackrel{\spadesuit}{\geq} \sup_{\varepsilon \in [0, 1]} \sum_{t=13}^T \mathbb{E} \left[2 \left(\frac{1}{2} - F_\varepsilon(\tilde{\alpha}_t(V_{\varepsilon, 1}, \dots, V_{\varepsilon, 2(t-1)})) \right)^2 \right] \\
 655 & \\
 656 & \stackrel{\circ}{\geq} \sum_{t=13}^T \mathbb{E} \left[2 \left(\frac{1}{2} - F_E(\tilde{\alpha}_t(V_{E, 1}, \dots, V_{E, 2(t-1)})) \right)^2 \right] \stackrel{\clubsuit}{\geq} 2 \sum_{t=13}^T \mathbb{E} \left[\left(\frac{5-2E}{8} - \tilde{\alpha}_t(V_{E, 1}, \dots, V_{E, 2(t-1)}) \right)^2 \right] \\
 657 & \\
 658 & \stackrel{\heartsuit}{\geq} 2 \sum_{t=13}^T \mathbb{E} \left[\left(\frac{5-2E}{8} - \mathbb{E} \left[\frac{5-2E}{8} \mid B_{E, 1}, \dots, B_{E, 2(t-1)}, D_1, \dots, D_{2(t-1)}, U_1, \dots, U_{2(t-1)} \right] \right)^2 \right] \\
 659 & \\
 660 & \stackrel{\spadesuit}{=} 2 \sum_{t=13}^T \mathbb{E} \left[\left(\frac{5-2E}{8} - \mathbb{E} \left[\frac{5-2E}{8} \mid B_{E, 1}, \dots, B_{E, 2(t-1)} \right] \right)^2 \right] = \frac{1}{8} \sum_{t=13}^T \mathbb{E} \left[(E - \mathbb{E}[E \mid B_{E, 1}, \dots, B_{E, 2(t-1)}])^2 \right] \\
 661 & \\
 662 & \stackrel{*}{=} \frac{1}{8} \sum_{t=13}^T \mathbb{E} \left[\left(E - \frac{\sum_{s=1}^{2(t-1)} B_{E, s} + 1}{2t} \right)^2 \right] = \frac{1}{8} \sum_{t=13}^T \int_0^1 \mathbb{E} \left[\left(\varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s} + 1}{2t} \right)^2 \right] d\varepsilon \\
 663 & \\
 664 & = \frac{1}{8} \sum_{t=13}^T \int_0^1 \mathbb{E} \left[\left(\varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} + \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2t} - \frac{1}{2t} \right)^2 \right] d\varepsilon \\
 665 & \\
 666 & = \frac{1}{8} \sum_{t=13}^T \int_0^1 \mathbb{E} \left[\left(\varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} + \frac{1}{2t(t-1)} \sum_{s=1}^{2(t-1)} B_{\varepsilon, s} - \frac{1}{2t} \right)^2 \right] d\varepsilon \\
 667 & \\
 668 & \geq \frac{1}{8} \sum_{t=13}^T \int_0^1 \mathbb{E} \left[\left(\varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} \right)^2 - 2 \left| \varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} \right| \left| \frac{1}{2t(t-1)} \sum_{s=1}^{2(t-1)} B_{\varepsilon, s} - \frac{1}{2t} \right| \right] d\varepsilon \\
 669 & \\
 670 & \geq \frac{1}{8} \sum_{t=13}^T \int_0^1 \mathbb{E} \left[\left(\varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} \right)^2 - \frac{1}{t} \left| \varepsilon - \frac{\sum_{s=1}^{2(t-1)} B_{\varepsilon, s}}{2(t-1)} \right| \right] d\varepsilon \\
 671 & \\
 672 & \stackrel{*}{\geq} \frac{1}{8} \sum_{t=13}^T \int_0^1 \left(\frac{\text{Var}(B_{\varepsilon, 1})}{2(t-1)} - \frac{1}{t} \sqrt{\frac{\text{Var}(B_{\varepsilon, 1})}{2(t-1)}} \right) d\varepsilon = \frac{1}{8} \sum_{t=13}^T \int_0^1 \left(\frac{\varepsilon(1-\varepsilon)}{2(t-1)} - \frac{1}{t} \sqrt{\frac{\varepsilon(1-\varepsilon)}{2(t-1)}} \right) d\varepsilon \\
 673 & \\
 674 & = \frac{1}{8} \sum_{t=13}^T \left(\frac{1}{12(t-1)} - \frac{\pi}{8t\sqrt{2(t-1)}} \right) \geq \frac{1}{8} \left(\frac{1}{12} - \frac{\pi}{16\sqrt{6}} \right) \sum_{t=12}^{T-1} \frac{1}{t} \geq \frac{1}{8} \left(\frac{1}{12} - \frac{\pi}{16\sqrt{6}} \right) \int_{12}^T \frac{1}{s} ds \\
 675 & \\
 676 & = \frac{1}{8} \left(\frac{1}{12} - \frac{\pi}{16\sqrt{6}} \right) \ln \left(\frac{T}{12} \right) \geq \frac{1}{16} \left(\frac{1}{12} - \frac{\pi}{16\sqrt{6}} \right) \ln(T). \\
 677 & \\
 678 & \\
 679 & \\
 680 & \\
 681 & \\
 682 & \\
 683 & \\
 684 &
 \end{aligned}$$

685 where “ \spadesuit ” follows from Lemma 1; “ \circ ” follows from the fact that E and $V_{\varepsilon, 1}, \dots, V_{\varepsilon, 2(t-1)}$ are inde-
 686 pendent of each other; “ \clubsuit ” follows from the fact that $\tilde{\alpha}_t$ takes values in $\left[\frac{1}{8}, \frac{7}{8}\right]$ and the explicit formula
 687 of F_ε in that interval for any $\varepsilon \in [0, 1]$; “ \heartsuit ” follows from the fact that $\tilde{\alpha}_t(V_{E, 1}, \dots, V_{E, 2(t-1)})$ is
 688 $\mathcal{F}_t := \sigma(B_{E, 1}, \dots, B_{E, 2(t-1)}, D_1, \dots, D_{2(t-1)}, U_1, \dots, U_{2(t-1)})$ -measurable and that, for any Y the
 689 minimizer in $L^2(\mathcal{F}_t)$ of the functional $X \mapsto \mathbb{E}[(Y - X)^2]$ is $X = \mathbb{E}[Y \mid \mathcal{F}_t]$; “ \diamond ” follows from the
 690 fact that E and $(D_1, \dots, D_{2(t-1)}, U_1, \dots, U_{2(t-1)})$ are independent of each other; “ $*$ ” follows from the
 691 fact $E \mid B_{E, 1}, \dots, B_{E, 2(t-1)}$ has a beta distribution; and “ $*$ ” follows from the fact that $B_{\varepsilon, 1}, B_{\varepsilon, 2}, \dots$
 692 is an i.i.d. Bernoulli process of parameter ε , together with Jensen’s inequality.
 693

694 B PROOF OF THEOREM 3

695 Without loss of generality, we can (and do!) assume that $T \geq 2$, and so $\log_2(MT) \geq 1$. First, let
 696 τ_T be the final value of τ if the algorithm ends at time T without a break, or define it as $\tau - 1$ if
 697 it ends with a break. For each $\tau \in [\tau_T]$, we define the epoch τ as the collection of rounds from $t_{\tau-1} + 1$
 698 to t_τ . Notice that, for each $\tau \in [\tau_T]$, we have that s_τ is the number of bits collected during the
 699 epoch τ . Let $Q_1^* := 1/2$, and define by induction $Q_{\tau+1}^* \leftarrow Q_\tau^* + \frac{1}{2^{\tau+1}}$ if $F(Q_\tau^*) < 1/2$, as $Q_\tau^* - \frac{1}{2^{\tau+1}}$
 700 if $F(Q_\tau^*) > 1/2$, or as Q_τ^* if $F(Q_\tau^*) = 1/2$. If there is $\tau \in \mathbb{N}$ such that $F(Q_\tau^*) = 1/2$, let $m := Q_\tau^*$.
 701

Otherwise, let $m \in [0, 1]$ be such that $F(m) = 1/2$ (its existence has already been pointed out after Lemma 1). Crucially, notice that for each $\tau \in \mathbb{N}$, we have that $|m - Q_\tau^*| \leq 2^{-\tau}$.

Let $(V_{x,k})_{x \in [0,1], k \in \mathbb{N}}$ be an independent family of random variables with common distribution given by ν , and for each $x \in [0, 1]$ and $t \in \mathbb{N}$, define $N_t(x) := 2 \cdot \sum_{k=1}^{t-1} \mathbb{I}\{P_k = x\}$. Notice that without loss of generality, we can assume that for each $t \in \mathbb{N}$ it holds that $V_{2t-1} := V_{P_t, N_t(P_t)+1}$ and $V_{2t} := V_{P_t, N_t(P_t)+2}$. Define the “good” event

$$\mathcal{E} := \bigcap_{i=1}^T \bigcap_{\substack{j=1 \\ j \text{ even}}}^T \left\{ \left| \frac{1}{j} \sum_{k=1}^j \mathbb{I}\{V_{Q_i^*, k} \leq Q_i^*\} - F(Q_i^*) \right| < \sqrt{\frac{\ln(2/\delta)}{2j}} \right\},$$

and notice that by De Morgan’s laws, a union bound, and Hoeffding’s inequality, we can upper bound the probability of the “bad” event \mathcal{E}^c by $\mathbb{P}[\mathcal{E}^c] \leq \delta T^2$. Notice that for each $i, j \in [T]$ with $F(Q_i^*) \neq \frac{1}{2}$ and j even satisfying $j \geq \frac{2 \ln(2/\delta)}{(\frac{1}{2} - F(Q_i^*))^2}$, then, whenever we are in the good event \mathcal{E} , we have that

$$\frac{1}{j} \sum_{k=1}^j \mathbb{I}\{V_{Q_i^*, k} \leq Q_i^*\} + \sqrt{\frac{\ln(2/\delta)}{2j}} < F(Q_i^*) + \sqrt{\frac{2 \ln(2/\delta)}{j}} \leq \frac{1}{2},$$

whenever $F(Q_i^*) < 1/2$, while

$$\frac{1}{j} \sum_{k=1}^j \mathbb{I}\{V_{Q_i^*, k} \leq Q_i^*\} - \sqrt{\frac{\ln(2/\delta)}{2j}} > F(Q_i^*) - \sqrt{\frac{2 \ln(2/\delta)}{j}} \geq \frac{1}{2}.$$

whenever $F(Q_i^*) > 1/2$. Instead, if $i, j \in [T]$ with $F(Q_i^*) = \frac{1}{2}$ and j is even, we have that

$$\frac{1}{j} \sum_{k=1}^j \mathbb{I}\{V_{Q_i^*, k} \leq Q_i^*\} + \sqrt{\frac{\ln(2/\delta)}{2j}} \geq F(Q_i^*) = \frac{1}{2}$$

and analogously

$$\frac{1}{j} \sum_{k=1}^j \mathbb{I}\{V_{Q_i^*, k} \leq Q_i^*\} - \sqrt{\frac{\ln(2/\delta)}{2j}} \leq F(Q_i^*) = \frac{1}{2}.$$

In particular, if we are in the good event \mathcal{E} , these inequalities imply on the one hand that $Q_1 = Q_1^*, \dots, Q_{\tau_T} = Q_{\tau_T}^*$ and, if $\tau \in [\tau_T]$ is such that $F(Q_\tau^*) = 1/2$, then $\tau = \tau_T$. On the other hand, if $\tau \in [\tau_T]$ is such that $F(Q_\tau^*) \neq 1/2$ and we are in the good event \mathcal{E} , they imply that the number of bits s_τ collected during the epoch τ cannot be greater than $\frac{2 \ln(2/\delta)}{(\frac{1}{2} - F(Q_\tau^*))^2}$, because the condition that ends the epoch τ with a break is met by the time that we have collected $\frac{2 \ln(2/\delta)}{(\frac{1}{2} - F(Q_\tau^*))^2}$ bits in that epoch.

Define $\tau_T^\# := \lceil \log_2(MT) \rceil$, define τ_T^b as the smallest $\tau \in \mathbb{N}$ such that $F(Q_\tau^*) = 1/2$ if it exists, and $+\infty$ otherwise, and define $\tau_T^* := \min(\tau_T^\#, \tau_T^b, \tau_T)$. In what follows, when we are in the event $\tau_T^\# > \max(\tau_T^b, \tau_T)$, we use the convention that any summation of the form $\sum_{\tau=\tau_T^*+1}^{\tau_T}$ is zero by definition. For each $t \in [T]$, define $\mathcal{H}_t := \sigma(V_1, \dots, V_{2t-2})$ as the σ -algebra generated by the history

756 observed before time t . We can control the regret in the following way
 757

$$\begin{aligned}
 758 \quad R_T &= \sum_{t=1}^T \mathbb{E}[G_t(m) - G_t(P_t)] = \sum_{t=1}^T \mathbb{E}\left[\mathbb{E}[G_t(m) - G_t(P_t) \mid \mathcal{H}_t]\right] \\
 759 & \\
 760 & \stackrel{\spadesuit}{=} \sum_{t=1}^T \mathbb{E}\left[\mathbb{E}[G_t(m) - G_t(p)]_{p=P_t}\right] \stackrel{\clubsuit}{=} 2 \cdot \sum_{t=1}^T \mathbb{E}\left[\left(\frac{1}{2} - F(P_t)\right)^2\right] \\
 761 & \\
 762 & \leq 2 \cdot \sum_{t=1}^T \mathbb{E}\left[\left(\frac{1}{2} - F(P_t)\right)^2 \mathbb{I}_{\mathcal{E}}\right] + \frac{T}{2} \cdot \mathbb{P}[\mathcal{E}^c] \\
 763 & \\
 764 & = \mathbb{E}\left[\sum_{\tau=1}^{\tau_T^* - 1} s_\tau \cdot \left(\frac{1}{2} - F(Q_\tau^*)\right)^2 \mathbb{I}_{\mathcal{E}}\right] + \mathbb{E}\left[\sum_{\tau=\tau_T^*}^{\tau_T} s_\tau \cdot \left(\frac{1}{2} - F(Q_\tau^*)\right)^2 \mathbb{I}_{\mathcal{E}}\right] + \frac{T}{2} \cdot \mathbb{P}[\mathcal{E}^c] \\
 765 & \\
 766 & \stackrel{\heartsuit}{\leq} \mathbb{E}\left[\sum_{\tau=1}^{\tau_T^* - 1} \frac{2 \ln(2/\delta)}{\left(\frac{1}{2} - F(Q_\tau^*)\right)^2} \cdot \left(\frac{1}{2} - F(Q_\tau^*)\right)^2 \mathbb{I}_{\mathcal{E}}\right] + \mathbb{E}\left[\sum_{\tau=\tau_T^*}^{\tau_T} s_\tau \cdot M^2 \cdot |m - Q_\tau^*|^2\right] + \frac{T}{2} \cdot \mathbb{P}[\mathcal{E}^c] \\
 767 & \\
 768 & \leq (\tau_T^\# - 1) \cdot 2 \cdot \ln(2/\delta) + T \cdot M^2 \cdot 2^{-2\tau_T^\#} + \delta \cdot \frac{T^3}{2} \leq 2 + 6 \log_2(MT) \ln(T),
 \end{aligned}$$

774 where in \spadesuit we used the Freezing Lemma (see, e.g., (Cesari & Colomboni, 2021, Lemma 8)), in \clubsuit we
 775 used Lemma 1, and in \heartsuit we used that fact that $F(m) = 1/2$ and F is M -Lipschitz.
 776

777 C PROOF OF THEOREM 4

779 We already know that algorithms that have access to full-feedback have to suffer worst-case regret
 780 of at least $c_1 \ln T$ if $T \geq c_2$, where c_1 and c_2 are the constants in the statement of Theorem 2. In
 781 particular, the same statement holds *a fortiori* for any 2-bit feedback algorithm, given that any 2-bit
 782 feedback algorithm can be trivially converted into an algorithm operating with full-feedback. It
 783 follows that it is enough to prove that there exist two universal constants \tilde{c}_1 and \tilde{c}_2 such that the worst-
 784 case regret of any 2-bit feedback algorithm is at least $\tilde{c}_1 \ln M$ whenever $T \geq \tilde{c}_2 \log_2(M)$. In fact, in
 785 this case, we can set $\bar{c}_1 := \frac{1}{2} \min(c_1, \tilde{c}_1)$ and $\bar{c}_2 := \max(c_2, \tilde{c}_2)$ to obtain that the worst-case regret of
 786 any 2-bit feedback algorithm is at least $2\bar{c}_1 \max(\ln T, \ln M) \geq \bar{c}_1 \ln(MT)$ whenever $T \geq \bar{c}_2 \log_2 M$.

787 We now prove the existence of \tilde{c}_1 and \tilde{c}_2 . Let $n \in \mathbb{N}$ be the greatest integer such that $2^n \leq M$ and
 788 consider the elements $\nu_k \in \mathcal{D}_M$ whose density is $2^n \cdot \mathbb{I}_{(\frac{k-1}{2^n}, \frac{k}{2^n})}$ for some $k \in [2^n]$, and notice that the
 789 corresponding cdfs are M -Lipschitz.
 790

791 Consider the following surrogate game. The adversary secretly chooses $k^* \in [2^n]$. The player action
 792 space is $[2^n]$. The surrogate game ends the first time $t \in \mathbb{N}$ when the player plays $I_t = k^*$. Before
 793 that, if the player plays $I_t \neq k^*$, the player suffers a loss $1/2$ and receives $\mathbb{I}\{I_t \leq k^*\}$ as feedback.
 794 Now, note that we can convert any algorithm α for the 2-bit feedback problem into an algorithm $\tilde{\alpha}$
 795 for the surrogate game in the following way. For each $k \in [2^n - 1]$, define $J_k := [(k-1)2^{-n}, k2^{-n})$
 796 and $J_{2^n} := [(2^n - 1)2^{-n}, 1]$. Whenever the algorithm α plays $P_t \in J_k$, the algorithm $\tilde{\alpha}$ plays $I_t := k$
 797 and passes $(\mathbb{I}\{I_t \leq k^*\}, \mathbb{I}\{I_t \leq k^*\})$ to α , where k^* is the underlying instance of the surrogate game.
 798 Now, notice that we can map every instance $k^* \in [2^n]$ for the surrogate game into the instance
 799 $\nu_{k^*} \in \mathcal{D}_M$ of the original problem and that the regret of the algorithm α on the instance ν_{k^*} is greater
 800 than or equal to than the regret of the algorithm $\tilde{\alpha}$ on the instance k^* . It follows that a worst-case
 801 regret lower bound for the surrogate game is also a worst-case regret lower bound for the original
 802 problem.

803 Fix an algorithm α for the surrogate game. Given that the surrogate game is deterministic, without
 804 any loss of generality we can assume that α is deterministic. We say that $S \subset [2^n]$ is a discrete
 805 segment if S is of the form $\{k \in [2^n] \mid a \leq k \leq b\}$ for some $a, b \in [2^n]$ with $a \leq b$. We can prove the
 806 following property by induction on $t = 0, 1, \dots, n-1$: there is a discrete segment J_t with at least
 807 $2^{n-t} - 1$ elements such that, for each $k, k' \in S_t$, the algorithm has not won the game by the time t
 808 and receives the same feedback (and hence selects the same actions) if the underlying instance is k
 809 or k' . For $t = 0$ the property is true by setting $S_0 := [2^n]$. Assume that the property is true for some
 $t \in \{0, 1, \dots, n-2\}$. Assume that $a, b \in [2^n]$ with $a \leq b$ are such that $S_t = \{k \in [2^n] \mid a \leq k \leq b\}$,
 where S_t is a segment that enjoys the property. Now, if the algorithm plays $I_{t+1} \notin S_t$ we can

810 set $S_{t+1} := S_t$, and we see that the required properties hold trivially. Instead, if $I_{t+1} \in S_t$ we set
 811 $S_{t+1} := \{k \in [2^n] \mid I_t + 1 \leq k \leq b\}$ if $I_{t+1} < \frac{a+b}{2}$ and we set $S_{t+1} := \{k \in [2^n] \mid a \leq k \leq I_t - 1\}$
 812 if $I_{t+1} \geq \frac{a+b}{2}$. Notice that given that S_t has at least $2^{n-t} - 1$ points, we have that S_{t+1} contains at
 813 least $\frac{2^{n-t}-2}{2} = 2^{n-(t+1)} - 1$ points and, for each $k \in S_{t+1}$, the game does not end by the time $t + 1$.
 814 Hence the induction step is proved. It follows that S_{n-1} is non-empty and, if we pick $k^* \in S_{n-1}$,
 815 the game goes on at least up to time $n - 1$ whenever the time horizon T is at least $n - 1$. Hence, if
 816 $T \geq \log_2(M)$ (which implies in particular that $T \geq n - 1$), the worst-case regret of the algorithm α is
 817 at least $\frac{n-1}{2} = \frac{n+1}{2} - 1 \geq \frac{\log_2(M)}{2} - 1 \geq \frac{\log_2(M)}{4} = \frac{1}{4 \ln(2)} \ln M$, where in the last inequality we used
 818 $M \geq 16$. Hence, we can pick $\tilde{c}_1 := \frac{1}{4 \ln 2}$ and $\tilde{c}_2 := 1$, concluding the proof.
 819
 820

821 D PROOF OF THEOREM 6

822
 823 For any $t \in \mathbb{N}$, tedious but straightforward computations show that
 824

$$825 \mathbb{P} \left[\sup_{p \in [0,1]} |\Psi(p) - \hat{\Psi}_t(p)| \geq \varepsilon \right] \leq \mathbb{P} \left[\sup_{p \in \mathbb{R}} \left| \frac{1}{2t} \sum_{s=1}^{2t} \mathbb{I}\{V_s \leq p\} - F(p) \right| \geq \frac{\varepsilon}{4} \right] \leq 2 \exp \left(-\frac{1}{4} \varepsilon^2 t \right),$$

826
 827 where the last inequality follows from the DKW inequality (Massart, 1990). Let $p^* \in$
 828 $\operatorname{argmax}_{p \in [0,1]} \Psi(p)$ (which does exist due to the upper-semicontinuity of Ψ). Then, for any $t \in \mathbb{N}$,
 829 we have that
 830

$$831 \mathbb{E} [\Psi(p^*) - \Psi(P_{t+1})] = \mathbb{E} [\Psi(p^*) - \hat{\Psi}_t(p^*)] + \underbrace{\mathbb{E} [\hat{\Psi}_t(p^*) - \hat{\Psi}_t(P_{t+1})]}_{\leq 0} + \mathbb{E} [\hat{\Psi}_t(P_{t+1}) - \Psi(P_{t+1})]$$

$$832$$

$$833 \leq 2 \mathbb{E} \left[\sup_{p \in [0,1]} |\Psi(p) - \hat{\Psi}_t(p)| \right] = 2 \int_0^{+\infty} \mathbb{P} \left[\sup_{p \in [0,1]} |\Psi(p) - \hat{\Psi}_t(p)| \geq \varepsilon \right] d\varepsilon$$

$$834$$

$$835 \leq 2 \int_0^{+\infty} 2 \exp \left(-\frac{1}{4} \varepsilon^2 t \right) d\varepsilon = \frac{4\sqrt{\pi}}{\sqrt{t}}.$$

$$836$$

$$837$$

$$838$$

839 Hence

$$840 R_T \leq 1 + \mathbb{E} \left[\sum_{t=2}^T (\Psi(p^*) - \Psi(P_t)) \right] \leq 1 + 4\sqrt{\pi} \sum_{t=1}^{T-1} \frac{1}{\sqrt{t}} \leq 1 + 8\sqrt{\pi} \cdot \sqrt{T-1}.$$

$$841$$

$$842$$

$$843$$

$$844$$

$$845$$

$$846$$

$$847$$

$$848$$

$$849$$

$$850$$

$$851$$

$$852$$

$$853$$

$$854$$

$$855$$

$$856$$

$$857$$

$$858$$

$$859$$

$$860$$

$$861$$

$$862$$

$$863$$