



A Contextual Online Learning Theory of Brokerage

François Bachoc, Tommaso R Cesari, Roberto Colombo

► To cite this version:

François Bachoc, Tommaso R Cesari, Roberto Colombo. A Contextual Online Learning Theory of Brokerage. 2024. hal-04619452

HAL Id: hal-04619452

<https://hal.science/hal-04619452>

Preprint submitted on 20 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Contextual Online Learning Theory of Brokerage

François Bachoc

Institut de Mathématiques de Toulouse, Université Paul Sabatier, Toulouse, France

Tommaso Cesari

University of Ottawa, Ottawa, Canada

Roberto Colomboni

Politecnico di Milano & Università degli Studi di Milano, Milano, Italy

June 20, 2024

Abstract

We study the role of *contextual information* in the online learning problem of brokerage between traders. At each round, two traders arrive with secret valuations about an asset they wish to trade. The broker suggests a trading price based on contextual data about the asset. Then, the traders decide to buy or sell depending on whether their valuations are higher or lower than the brokerage price.

We assume the market value of traded assets is an unknown linear function of a d -dimensional vector representing the contextual information available to the broker. Additionally, we model traders' valuations as independent bounded zero-mean perturbations of the asset's market value, allowing for potentially different unknown distributions across traders and time steps. Consistently with the existing online learning literature, we evaluate the performance of a learning algorithm with the regret with respect to the *gain from trade*. If the noise distributions admit densities bounded by some constant L , then, for any time horizon T :

- If the agents' valuations are revealed after each interaction, we provide an algorithm achieving $O(Ld\ln T)$ regret, and show a corresponding matching lower bound of $\Omega(Ld\ln T)$.
- If only their willingness to sell or buy at the proposed price is revealed after each interaction, we provide an algorithm achieving $O(\sqrt{LdT \ln T})$ regret, and show that this rate is optimal (up to logarithmic factors), via a lower bound of $\Omega(\sqrt{LdT})$.

To complete the picture, we show that if the bounded density assumption is lifted, then the problem becomes unlearnable, even with full feedback.

1 Introduction

Inspired by a recent stream of literature [3, 7, 9, 12, 14, 15], we approach the bilateral trade problem of brokerage between traders through the lens of online learning. When viewed from a regret minimization perspective, bilateral trade has been explored over rounds of seller/buyer interactions with no prior knowledge of their private valuations. As in [9], we focus on the case where traders are willing to either buy or sell, depending on whether their valuations for the asset being traded are above or below the brokerage price.

This setting is especially relevant for over-the-counter (OTC) markets. Serving as alternatives to conventional exchanges, OTC markets operate in a decentralized manner and are a vital part of the global financial landscape.¹ In contrast to centralized exchanges, the lack of strict protocols and regulations allows brokers to take on the responsibility of bridging the gap between buyers and sellers, who may not have direct access to one another. In addition to facilitating interactions between parties, brokers leverage their contextual

¹In the US alone, the value of assets traded in OTC markets exceeded a remarkable 50 trillion USD in 2020, surpassing centralized markets by more than 20 trillion USD [54]. This growth has been steadily increasing since 2016 [1].

knowledge and market insights to determine appropriate pricing for assets. By examining factors such as supply and demand, market trends, and other asset-specific information, brokers aim to propose prices that reflect the true value of the asset being traded. This price discovery process is a crucial aspect of a broker's role, as it helps ensure efficient transactions by accounting for the unique circumstances surrounding each asset. Additionally, in many OTC markets, as in our setting, traders choose to either buy or sell depending on the contingent market conditions [51]. This behavior is observed across a broad range of asset trades, including stocks, derivatives, art, collectibles, precious metals and minerals, energy commodities like gas and oil, and digital currencies (cryptocurrencies), among others [9].

In the existing literature on online learning for bilateral trade, the contextual version of this problem has never been investigated. This case is of significant interest given that the broker often has access to meaningful information about the asset being traded and the surrounding market conditions *before* having to propose a trading price. This information might help the broker to propose more targeted trading prices by inferring the current market value of the corresponding asset, and ignoring it could be extremely costly in terms of missing trading opportunities. We aim to fill this gap in the online learning literature on bilateral trade to guide brokers in these contextual scenarios.

1.1 Setting

In the following, the elements of any Euclidean space are treated as column vectors and, for any real number x, y , we denote their minimum by $x \wedge y$ and their maximum by $x \vee y$.

We study the following problem. At each time $t \in \mathbb{N}$,

- Two traders arrive with private valuations $V_t, W_t \in [0, 1]$ about an asset they want to trade.
- The broker observes a context $c_t \in [0, 1]^d$ and proposes a trading price $P_t \in [0, 1]$.
- If the price P_t lies between the lowest valuation $V_t \wedge W_t$ and highest valuation $V_t \vee W_t$ (meaning the trader with the minimum valuation is ready to sell at P_t and the trader with the maximum valuation is eager to buy at P_t), the asset is bought by the trader with the highest valuation from the trader with the lowest valuation at the brokerage price P_t .
- Some feedback is disclosed.

We assume that at each time $t \in \mathbb{N}$, the traders' valuations V_t and W_t are zero-mean perturbations of some (unknown) market value $m_t \in [0, 1]$ of the asset being traded. Specifically, we assume that there exists an independent family of random variables $(\xi_t, \zeta_t)_{t \in \mathbb{N}}$ such that, for each $t \in \mathbb{N}$, it holds that $\mathbb{E}[\xi_t] = 0 = \mathbb{E}[\zeta_t]$ and $V_t = m_t + \xi_t$ and $W_t = m_t + \zeta_t$.² In particular, we note that $\mathbb{E}[V_t] = m_t = \mathbb{E}[W_t]$, for each $t \in \mathbb{N}$, and hence the traders' valuations can be thought of as zero-mean perturbations of the market value of the asset due to individual preferences.

Following the recent stream of bilateral trade literature investigating the interplay between learning and the regularity of the underlying valuation distributions [9, 14, 15], we focus on the case when the traders' valuation distributions admit densities that are uniformly bounded by some constant $L \geq 1$. We note that this assumption is equivalent to the same uniformly bounded density assumption on the distributions of the noise $\xi_1, \zeta_1, \xi_2, \zeta_2, \dots$. Then, we analyze what happens when the bounded density assumption is lifted.

We assume an unknown linear relation exists between the market value m_t for the asset being traded at time t and the corresponding context c_t the broker observes before proposing a trading price. Specifically, we assume that there exists $\phi \in [0, 1]^d$, unknown to the broker, such that, for each $t \in \mathbb{N}$, it holds that $m_t = c_t^\top \phi$.

We model the sequence of contexts $(c_t)_{t \in \mathbb{N}}$ as a deterministic $[0, 1]^d$ -valued sequence, possibly generated in an adversarial manner by someone who knows the broker's algorithm. As explained above, we assume that the sequence of contexts is initially unknown but sequentially discovered by the broker.

²We remark that we are not assuming that the two processes $(\xi_t)_{t \in \mathbb{N}}$ and $(\zeta_t)_{t \in \mathbb{N}}$ are i.i.d., and in fact the distributions of these random variables may change over time.

Consistently with the existing bilateral trade literature, the reward associated with each interaction is the sum of the net utilities of the traders, known as *gain from trade*. Formally, for any $p, v, w \in [0, 1]$, the utility of a price p when the valuations of the traders are v and w is

$$g(p, v, w) := (\underbrace{v \vee w - p}_{\text{buyer's net gain}} + \underbrace{p - v \wedge w}_{\text{seller's net gain}}) \mathbb{I}\{\underbrace{v \wedge w \leq p \leq v \vee w}_{\text{whenever a trade happens}}\} = (v \vee w - v \wedge w) \mathbb{I}\{v \wedge w \leq p \leq v \vee w\}.$$

The aim of the learner is to minimize the *regret* with respect to the best function of the contexts, defined, for any time horizon $T \in \mathbb{N}$, as

$$R_T := \sup_{p^* : [0, 1]^d \rightarrow [0, 1]} \mathbb{E}\left[\sum_{t=1}^T \left(\text{GFT}_t(p^*(c_t)) - \text{GFT}_t(P_t)\right)\right],$$

where we let $\text{GFT}_t(p) := g(p, V_t, W_t)$ for all $p \in [0, 1]$, and the expectation is taken with respect to the randomness in $(\xi_t, \zeta_t)_{t \in \mathbb{N}}$ and, possibly, the internal randomization used to choose the trading prices $(P_t)_{t \in \mathbb{N}}$.

Finally, we consider the two most studied types of feedback in the bilateral trade literature. Specifically, at each round t , only after having posted the price P_t , the learner receives either:

- *Full feedback*, i.e., the valuations V_t and W_t of the two current traders are disclosed.
- *Two-bit feedback*, i.e., only the indicator functions $\mathbb{I}\{P_t \leq V_t\}$ and $\mathbb{I}\{P_t \leq W_t\}$ are disclosed.

The information gathered in the full feedback model reflects *direct revelation mechanisms*, where traders disclose their valuations V_t and W_t prior to each round, but the price determined by the mechanism at time t is based solely on the previous valuations $V_1, W_1, \dots, V_{t-1}, W_{t-1}$. Conversely, the two-bit feedback model reflects *posted price* mechanisms. In this model, traders only indicate their willingness to buy or sell at the posted price, and their valuations V_t and W_t remain undisclosed.

1.2 Our contributions

Under the assumption that the traders' valuations are unknown linear functions of d -dimensional contexts perturbed by zero-mean noise with time-variable densities bounded by some L , and with the goal of designing *simple* and *interpretable* optimal algorithms, we make the following contributions.

1. We prove a structural result (Lemma 1) with two crucial consequences. First, Lemma 1 shows that posting the traders' (unknown) expected valuation as the trading price would maximize the expected gain from trade. Second, it proves that the loss paid by posting a suboptimal price is at most quadratic in the distance from an optimal one.
2. In the full feedback setting, we introduce an algorithm based on ridge regression estimation (Algorithm 1) and, leveraging the previous lemma, we prove its optimality by showing matching $Ld \ln T$ regret upper and lower bounds (Theorems 1 and 2).
3. In the two-bit feedback setting, the prices we post directly affect the information we retrieve. We note that this information is so scarce that it is not even enough to reconstruct bandit feedback. We solve this challenging exploration-exploitation dilemma by proposing an algorithm (Algorithm 2) that decides to either explore or exploit adaptively, based on the amount of contextual information gathered so far, and prove its optimality by showing a $\sqrt{LdT \ln T}$ regret upper bound (Theorem 3) and a matching (up to a $\sqrt{\ln T}$) \sqrt{LdT} lower bound (Theorem 4).
4. Finally, we investigate the necessity of the bounded density assumption: by lifting this assumption, we show that the problem becomes unlearnable (Theorem 5).

To the best of our knowledge, our work is the first to analyze a *non-parametric noisy contextual* bilateral trade problem (in fact, the first that analyzes a *contextual* bilateral trade problem in general) and one of only two works on bilateral trade (the other one being [9]) where the dependence on *all* relevant parameters is tight. As we discuss in Section 1.3, most related works on non-contextual bilateral trade obtain (at best) a matching dependence in the time horizon only, while those on non-parametric noisy contextual pricing/auctions lack matching lower bounds altogether.

1.3 Related Works

Building upon the foundational work of Myerson and Satterthwaite [47], a rich body of research has investigated bilateral trade from a game-theoretic and best-approximation standpoint [2, 4, 8, 10, 20, 21, 22, 23, 24, 32]. For an insightful analysis of this literature, see [12].

Our work builds upon the recent research on bilateral trade within online learning settings. Given the close relationship between our and these existing works, we discuss these connections in detail. First, to the best of our knowledge, the existing online learning literature on bilateral trade *never* discussed contextual problems. In [3, 7, 12, 13, 14, 15], the authors studied non-contextual bilateral trade problems where sellers and buyers have definite roles. [12, 14] show that the adversarial setting is unlearnable, and hence they focus on the case where sellers' and buyers' valuations form an i.i.d. process. They obtain a \sqrt{T} regret rate in the full-feedback setting. For the two-bit feedback case, they show that the problem is unlearnable in general, but it turns out to be learnable at a tight regret rate of $T^{2/3}$ by assuming that sellers' and buyers' valuations are independent of each other and they admit a uniformly bounded density. [3] show that learning is achievable in the adversarial case if the weaker α -regret objective is considered. Specifically, in the full-feedback case, they obtain a tight 2-regret rate of \sqrt{T} . In the two-bit feedback case, they show that learning is impossible in general, but by allowing the learner to use weakly budget-balanced mechanisms, they recover a 2-regret of order $T^{3/4}$, without a matching lower bound. In a different direction, [13, 15] show that learning is achievable in the adversarial case if the adversary is forced to be *smooth*, i.e., the sellers' and buyers' valuation distributions may change adversarially over time, but these distributions admit uniformly bounded densities. In the full-feedback case, they obtain a tight \sqrt{T} regret rate. In the two-bit feedback case, they show that the problem is still unlearnable, but, by allowing the learner to use weakly budget-balanced mechanisms, they prove a surprisingly sharp $T^{3/4}$ regret rate. [7] propose the notion of globally budget-balanced mechanisms, a further relaxation of the weakly budget-balanced notion, under which they show that learning is achievable in the adversarial case at a tight regret rate of \sqrt{T} in the full-feedback case, and at a regret rate of $T^{3/4}$ in the two-bit feedback case, without a matching lower bound. We remark that in all the papers we discussed so far, every two-bit feedback upper bound that requires a bounded density assumption lacks a corresponding lower bound with a sharp dependence on this parameter. The closest to our setting is the one proposed in [9]. There, the authors study the non-contextual version of our trading problem with flexible sellers' and buyers' roles, with the further assumption that the sellers' and buyers' valuations form an i.i.d. sequence. Under the M -bounded density assumption, they obtain tight $M \ln T$ and \sqrt{MT} regret rates in the full-feedback and two-bit feedback settings, respectively. If the bounded density assumption is removed, they show that the learning rate degrades to \sqrt{T} in the full-feedback case and the problem turns out to be unlearnable in the two-bit feedback case. We remark that, interestingly, under the bounded density assumption, we are able to achieve the same regret rates in the contextual version of this problem without requiring that traders share the same valuation distribution, while, without the bounded density assumption, the contextual problem is unlearnable even under full-feedback.

Our linear assumption appears commonly in the literature on digital markets, particularly in problems like pricing and auctions. In [18, 19], the authors first address a deterministic setting, then a noisy one with *known* noise distribution where they obtain a regret rate of order $T^{2/3}$ without presenting a lower bound. The deterministic case has also been investigated in [39, 40, 41, 42, 43].

Notably, the best results currently known only apply to deterministic settings, while, in the case of noisy linear functions, to the best of our knowledge [5, 11, 17, 25, 31, 45, 50, 56], the only known guarantees are limited to parametric or semi-parametric settings, a clear general picture of the minimax rates is still missing, and only slow upper bounds of order $T^{2/3}$ are known in the non-parametric case where the noise distribution

is unknown. In contrast, thanks to our Lemma 1, we are able to address the trading problem even in the non-parametric noisy linear case, obtaining optimal rates (matched by corresponding lower bounds) which are significantly faster than the ones known for pricing.

Another rich related field explored in its many variants [26, 27, 28, 29, 30, 33, 34, 35, 36, 38, 44, 46, 48, 49, 52, 53, 57] is contextual linear bandits. In its standard form, at the beginning of each round, an action set is revealed to the learner, and the assumption is that the reward (which equals the feedback) is a linear function of the action selected from the action set. Instead, in our setting, the market price is a linear function of the context, while the rewards are linked to the price the learner posts by the non-linear gain from trade function. Moreover, in contrast to contextual bandits, in our 2-bit feedback model, the feedback differs from and is not sufficient to compute the reward of the action the learner selects at every round. For these reasons, the techniques appearing in contextual linear bandits do not directly translate to our problem.

2 Structural Results

We begin by presenting a structural result whose economic interpretation is as follows: even if the broker does not know the traders' valuation distribution, if these valuations can be modeled as zero-mean noisy perturbations with bounded densities of some market value, then the best price to post to maximize the expected gain from trade is precisely the market value. In particular, this generalizes a similar result appearing in [9], which holds under the further assumption that the valuations have the exact same distribution. The following result also gives a representation formula for the expected gain from trade, which implies in particular that the cost of posting a suboptimal price is only quadratic in the distance from the market value. This structural result is the key to unraveling the intricacies of the noisy contextual setting, and it is what ultimately allows us to obtain tight regret guarantees in all settings, distinguishing ours from similar contextual pricing works.

Lemma 1. *Suppose that V and W are two $[0, 1]$ -valued independent random variables, with possibly different densities bounded by some constant $L \geq 1$, and such that $\mathbb{E}[V] = \mathbb{E}[W] = m$. Then, for each $p \in [0, 1]$, it holds that*

$$0 \leq \mathbb{E}[g(m, V, W) - g(p, V, W)] \leq L|m - p|^2.$$

Proof. We denote by F (resp., G) the cumulative distribution function of V (resp., W). For each $p \in [0, 1]$, from the Decomposition Lemma in [12, Lemma 1], it holds that

$$\begin{aligned} \mathbb{E}[(W - V)\mathbb{I}\{V \leq p \leq W\}] &= F(p) \int_p^1 (1 - G(\lambda)) d\lambda + (1 - G(p)) \int_0^p F(\lambda) d\lambda, \\ \mathbb{E}[(V - W)\mathbb{I}\{W \leq p \leq V\}] &= G(p) \int_p^1 (1 - F(\lambda)) d\lambda + (1 - F(p)) \int_0^p G(\lambda) d\lambda. \end{aligned}$$

Hence, for each $p \in [0, 1]$,

$$\begin{aligned} \mathbb{E}[(W - V)\mathbb{I}\{V \leq p \leq W\}] &= F(p) \int_p^1 (1 - G(\lambda)) d\lambda + (1 - G(p)) \int_0^p F(\lambda) d\lambda \\ &= F(p) \left(m - \int_0^p (1 - G(\lambda)) d\lambda \right) + \int_0^p F(\lambda) d\lambda - G(p) \int_0^p F(\lambda) d\lambda \\ &= \int_0^p F(\lambda) d\lambda + (m - p)F(p) - pG(p) + G(p) \int_0^p (1 - F(\lambda)) d\lambda + F(p) \int_0^p G(\lambda) d\lambda \\ &= \int_0^p (F + G)(\lambda) d\lambda + (m - p)(F + G)(p) - G(p) \left(m - \int_0^p (1 - F(\lambda)) d\lambda \right) + (F(p) - 1) \int_0^p G(\lambda) d\lambda \\ &= \int_0^p (F + G)(\lambda) d\lambda + (m - p)(F + G)(p) - \left(G(p) \int_p^1 (1 - F(\lambda)) d\lambda + (1 - F(p)) \int_0^p G(\lambda) d\lambda \right) \\ &= \int_0^p (F + G)(\lambda) d\lambda + (m - p)(F + G)(p) - \mathbb{E}[(V - W)\mathbb{I}\{W \leq p \leq V\}]. \end{aligned}$$

Rearranging, it follows that, for each $p \in [0, 1]$,

$$\begin{aligned}\mathbb{E}[g(p, V, W)] &= \mathbb{E}[(W - V)\mathbb{I}\{V \leq p \leq W\}] + \mathbb{E}[(V - W)\mathbb{I}\{W \leq p \leq V\}] \\ &= \int_0^p (F + G)(\lambda) d\lambda + (m - p)(F + G)(p).\end{aligned}$$

Hence, for any $p \in [0, 1]$, it holds that

$$\mathbb{E}[g(m, V, W) - g(p, V, W)] = \int_p^m ((F + G)(\lambda) - (F + G)(p)) d\lambda \geq 0.$$

Finally, since F and G are absolutely continuous with weak derivative bounded by L , by the fundamental theorem of calculus [6, Theorem 14.16] it holds that, for $p \in [0, 1]$,

$$\mathbb{E}[g(m, V, W) - g(p, V, W)] = \int_p^m \int_p^\lambda (F' + G')(\vartheta) d\vartheta d\lambda \leq 2L \int_p^m |\lambda - p| d\lambda = L|m - p|^2. \quad \square$$

As a corollary of Lemma 1, we obtain the following result, that upper bounds the regret in terms of the sum of the squared distances between the prices the algorithm posts and the actual market values.

Corollary 1. *Consider the setting introduced in Section 1.1. If the valuations admit densities bounded by a constant $L \geq 1$, then, for any time horizon $T \in \mathbb{N}$, we have*

$$R_T = \mathbb{E} \left[\sum_{t=1}^T (\text{GFT}_t(c_t^\top \phi) - \text{GFT}_t(P_t)) \right] \leq \sum_{t=1}^T 1 \wedge (L \mathbb{E} [|P_t - c_t^\top \phi|^2]).$$

Proof. Given that for each $t \in \mathbb{N}$ and each $p \in [0, 1]$ it holds that $\text{GFT}_t(p) \in [0, 1]$, we have

$$\sup_{p \in [0, 1]} \mathbb{E}[\text{GFT}_t(p) - \text{GFT}_t(P_t)] \leq 1,$$

and hence, recalling that $m_t = c_t^\top \phi$ and that $\mathbb{E}[V_t] = m_t = \mathbb{E}[W_t]$, we also have, for each $T \in \mathbb{N}$,

$$\begin{aligned}R_T &= \sup_{p^* : [0, 1]^d \rightarrow [0, 1]} \sum_{t=1}^T 1 \wedge (\mathbb{E}[g(p^*(c_t), V_t, W_t)] - \mathbb{E}[g(P_t, V_t, W_t)]) \\ &\stackrel{(\circ)}{=} \sum_{t=1}^T 1 \wedge (\mathbb{E}[g(c_t^\top \phi, V_t, W_t)] - \mathbb{E}[g(P_t, V_t, W_t)]) \\ &\stackrel{(*)}{=} \sum_{t=1}^T 1 \wedge \mathbb{E} \left[\left[\mathbb{E}[g(c_t^\top \phi, V_t, W_t) - g(p, V_t, W_t)] \right]_{p=P_t} \right] \stackrel{(\circ)}{\leq} \sum_{t=1}^T 1 \wedge (L \mathbb{E} [|P_t - c_t^\top \phi|^2]),\end{aligned}$$

where (\circ) follows from Lemma 1, and $(*)$ from the Freezing Lemma [16, Lemma 8]. \square

3 Full Feedback

In this section, we focus on the full feedback setting, corresponding to direct revelation mechanisms. We show that performing ridge regression to obtain an estimate of the unknown vector ϕ and using it as a proxy linear function to convert contexts into prices (Algorithm 1) is enough to achieve logarithmic regret. In the following, we denote by $\mathbf{1}_d$ the d -dimensional identity matrix.

Theorem 1. *Consider the full-feedback setting introduced in Section 1.1. If the learner runs Algorithm 1 and the traders' valuations admit a density bounded by $L \geq 1$, then, for any time horizon $T \in \mathbb{N}$, it holds that $R_T \leq 1 + 4Ld \ln T$.*

Algorithm 1: Ridge Regression Pricing — Full Feedback

Observe context c_1 , post $P_1 \doteq 1/2$, and receive feedback V_1, W_1 ;
Let $x_1 \doteq [c_1 | c_1]$, let $Y_1 \doteq [V_1 | W_1]$, and compute $\hat{\phi}_1 \doteq (x_1 x_1^\top + d^{-1} \mathbf{1}_d)^{-1} x_1 Y_1^\top$;
for time $t = 2, 3, \dots$ **do**
 Observe context c_t , post $P_t \doteq c_t^\top \hat{\phi}_{t-1}$, and receive feedback V_t, W_t ;
 Let $x_t \doteq [x_{t-1} | c_t | c_t]$, $Y_t \doteq [Y_{t-1} | V_t | W_t]$, and compute $\hat{\phi}_t \doteq (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} x_t Y_t^\top$;

Proof. Recall that $(\xi_t, \zeta_t)_{t \in \mathbb{N}}$ is an independent family of zero mean random variables each of them admitting a density bounded by L , that for any $t \in \mathbb{N}$, it holds that $m_t = c_t^\top \phi$, that $m_t + \xi_t = V_t \in [0, 1]$ and that $m_t + \zeta_t = W_t \in [0, 1]$. For any $t \in \mathbb{N}$, simple calculations show that

$$\mathbb{E}[|c_{t+1}^\top \hat{\phi}_t - c_{t+1}^\top \phi|^2] = \underbrace{\left(\mathbb{E}[c_{t+1}^\top \hat{\phi}_t - c_{t+1}^\top \phi] \right)^2}_{\text{bias}} + \underbrace{\text{Var}[c_{t+1}^\top \hat{\phi}_t]}_{\text{variance}}.$$

which is the well-known decomposition of the quadratic error with bias and variance of the estimator $c_{t+1}^\top \hat{\phi}_t$ for the quantity $c_{t+1}^\top \phi$. Noting that, for each $t \in \mathbb{N}$, it holds that $\mathbb{E}[Y_t^\top] = x_t^\top \phi$, we have,

$$\begin{aligned} \mathbb{E}[c_{t+1}^\top \hat{\phi}_t - c_{t+1}^\top \phi] &= c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} x_t x_t^\top \phi - c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} (x_t x_t^\top \phi + d^{-1} \phi) \\ &= -c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} d^{-1} \phi \doteq (\circ), \end{aligned}$$

and hence, by the Cauchy-Schwarz inequality applied to the scalar product $(a, b) \mapsto a^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} b$, by the fact that $(x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} \leq d^{-1} \mathbf{1}_d^{-1}$ (where, for any two symmetric matrices A_1, A_2 , we say that $A_1 \leq A_2$ if and only if $A_2 - A_1$ is semi-positive definite), and by the fact that $\|\phi\|_2^2 \leq d$, we can control the bias term as follows

$$\begin{aligned} (\mathbb{E}[c_{t+1}^\top \hat{\phi}_t - c_{t+1}^\top \phi])^2 &= (\circ)^2 \leq c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1} \cdot d^{-1} \phi^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} d^{-1} \phi \\ &\leq c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1} \cdot d^{-1} \phi^\top (d^{-1} \mathbf{1}_d)^{-1} d^{-1} \phi \leq c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1}. \end{aligned} \quad (1)$$

For each $t \in \mathbb{N}$, letting Δ_t be the $2t \times 2t$ diagonal matrix with vector of diagonal elements given by $(\text{Var}[V_1], \text{Var}[W_1], \text{Var}[V_2], \text{Var}[W_2], \dots, \text{Var}[V_t], \text{Var}[W_t])$, we have

$$\text{Var}[c_{t+1}^\top \hat{\phi}_t] = c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} (x_t \Delta_t x_t^\top) (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1}. \quad (2)$$

Now, for each $t \in \mathbb{N}$, given that $V_1, W_1, \dots, V_t, W_t$ are $[0, 1]$ -valued, we have that Δ_t is diagonal with diagonal elements less than 1, and hence $x_t \Delta_t x_t^\top \leq x_t x_t^\top + d^{-1} \mathbf{1}_d$, which yields a control on the variance term as follows,

$$\text{Var}[c_{t+1}^\top \hat{\phi}_t] \leq c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} (x_t x_t^\top + d^{-1} \mathbf{1}_d) (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1} = c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1}.$$

In the end, for each $t \in \mathbb{N}$, we have

$$\begin{aligned} \mathbb{E}[|c_{t+1}^\top \hat{\phi}_t - c_{t+1}^\top \phi|^2] &\leq 2c_{t+1}^\top (x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1} c_{t+1} = 2 \|c_{t+1}\|_{(x_t x_t^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 \\ &= 2 \|c_{t+1}\|_{(2 \sum_{s=1}^t c_s c_s^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 = \left\| \sqrt{2} c_{t+1} \right\|_{(\sum_{s=1}^t (\sqrt{2} c_s) (\sqrt{2} c_s)^\top + d^{-1} \mathbf{1}_d)^{-1}}^2, \end{aligned} \quad (3)$$

where, for any positive definite matrix $A \in \mathbb{R}^{d \times d}$ and each $u \in \mathbb{R}^d$, we have defined $\|u\|_A \doteq \sqrt{u^\top A u}$. Now, for any time horizon $T \in \mathbb{N}$, leveraging Corollary 1, we have that

$$\begin{aligned} R_T &\leq \sum_{t=1}^T 1 \wedge \left(L \mathbb{E}[|P_t - c_t^\top \phi|^2] \right) \leq 1 + \sum_{t=1}^{T-1} 1 \wedge \left(L \mathbb{E}[|c_{t+1}^\top \hat{\phi}_t - c_{t+1}^\top \phi|^2] \right) \\ &\leq 1 + L \sum_{t=1}^{T-1} 1 \wedge \left\| \sqrt{2} c_{t+1} \right\|_{(\sum_{s=1}^t (\sqrt{2} c_s) (\sqrt{2} c_s)^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 \doteq (\star). \end{aligned}$$

From here, we apply the elliptical potential lemma [37, Lemma 19.4] to obtain that, for any time horizon $T \in \mathbb{N}$,

$$R_T \leq (\star) \leq 1 + 2Ld \ln \left(\frac{dd^{-1} + 2d(T-1)}{dd^{-1}} \right) = 1 + 2Ld \ln(1 + 2d(T-1)) \leq 1 + 2Ld \ln(2dT).$$

If $d < T/2$, this implies that $R_T \leq 1 + 2Ld \ln(2dT) \leq 1 + 4Ld \ln T$. If, instead, $d \geq T/2$, then, recalling that $L \geq 1$, we obtain once again that $R_T \leq T \leq 1 + 4Ld \ln T$, concluding the proof. \square

We conclude this section by stating a matching worst-case $\Omega(Ld \ln T)$ regret lower bound for any algorithm in the full-feedback case, proving the optimality of Algorithm 1.

At a high level, the proof of this result is based on first building a sequence of contexts defined as a common element of the canonical basis of \mathbb{R}^d during each one of d blocks of T/d consecutive time-steps. Then, in each block, an adaptation of the non-contextual full-feedback lower bound construction in [9, Theorem 3] yields a lower bound of order $L \ln(T/d)$. Summing over blocks gives the result. For a full proof of this result, see Appendix A.

Theorem 2. *There exist two numerical constants $a, b > 0$ such that, for any $L \geq 2$ and any time horizon $T \geq \max(4, adL^5, 2d)$, there exists a sequence of contexts $c_1, \dots, c_T \in [0, 1]^d$ such that, for any algorithm α for the contextual brokerage problem with full feedback, there exists a vector $\phi \in [0, 1]^d$ and two zero-mean independent sequences $(\xi_t)_{t \in [T]}$ and $(\zeta_t)_{t \in [T]}$ independent of each other, such that if we define $V_t := c_t^\top \phi + \xi_t$ and $W_t := c_t^\top \phi + \zeta_t$, then for each $t \in [T]$ it holds that $c_t^\top \phi \in [0, 1]$, V_t and W_t are $[0, 1]$ -valued random variables with density bounded by L , and the regret of α on the sequence of traders' valuations $V_1, W_1, \dots, V_T, W_T$ satisfies $R_T \geq bLd \ln T$.*

We remark that the previous lower bound holds even for algorithms that have prior knowledge of the sequence of contexts c_1, c_2, \dots and that Theorem 1 shows that Algorithm 1 matches the optimal $Ld \ln T$ rate even without this *a-priori* knowledge.

4 Two-bit Feedback

In this section, we focus on the two-bit feedback setting, corresponding to posted-price mechanisms. We show that a simple deterministic rule that decides to either explore (by posting a price drawn uniformly in $[0, 1]$ to gather feedback to reconstruct the cumulative distribution functions of the traders' valuations) or exploit (by posting the scalar product of the context and the current ridge regression estimate of the unknown weight vector ϕ) based on the amount of information gathered along the various context dimensions (Algorithm 2) is enough to achieve $\tilde{\mathcal{O}}(\sqrt{LdT})$ regret. We recall that $\mathbf{1}_d$ is the d -dimensional identity matrix. Also, for any positive definite matrix $A \in \mathbb{R}^{d \times d}$, we define $\|\cdot\|_A : \mathbb{R}^d \rightarrow [0, \infty)$, $v \mapsto \sqrt{v^\top A v}$.

Algorithm 2: Scouting Ridge Regression Pricing — Two-bit Feedback

Post P_1 uniformly at random in $[0, 1]$, and observe $D_1 := \mathbb{I}\{P_1 \leq V_1\}$, $E_1 := \mathbb{I}\{P_1 \leq W_1\}$;

Let $b_1 := 1$, let $x_1 := [c_1 \mid c_1]$, let $Y_1 := [D_1 \mid E_1]$ and compute $\hat{\phi}_1 := (x_1 x_1^\top + d^{-1} \mathbf{1}_d)^{-1} x_1 Y_1^\top$;

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

Observe context c_t and define $b_t := \mathbb{I}\left\{ \|\sqrt{2}c_t\|_{(x_{t-1} x_{t-1}^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 > \sqrt{\frac{2d \ln(1+2d(T-1))}{LT}} \right\}$;

if $b_t = 1$ **then**

Post P_t uniformly at random in $[0, 1]$, and observe $D_t := \mathbb{I}\{P_t \leq V_t\}$, $E_t := \mathbb{I}\{P_t \leq W_t\}$;

Let $x_t := [x_{t-1} \mid c_t \mid c_t]$, let $Y_t := [Y_{t-1} \mid D_t \mid E_t]$ and compute $\hat{\phi}_t := (x_t x_t^\top + \mathbf{1}_d)^{-1} x_t Y_t^\top$;

else

Post $P_t = c_t^\top \hat{\phi}_{t-1}$ and let $x_t := x_{t-1}$, $Y_t := Y_{t-1}$, and $\hat{\phi}_t := \hat{\phi}_{t-1}$;

Theorem 3. Consider the two-bit feedback setting introduced in Section 1.1. If the learner runs Algorithm 2 and the traders' valuations admit a density bounded by $L \geq 1$, then, for any time horizon T such that $LT \geq 2d\ln(1 + 2d(T - 1))$, it holds that $R_T \leq 1 + 4\sqrt{LdT \ln T}$.

Proof. Without loss of generality we assume that $T \geq 2$. Note that for any $t \in \mathbb{N}$, if $b_t = 1$, then

$$\mathbb{E}[D_t] = \mathbb{P}[P_t \leq V_t] = \int_0^1 \mathbb{P}[u \leq V_t] du = \mathbb{E}[V_t] = \mathbb{E}[c_t^\top \phi + \xi_t] = c_t^\top \phi,$$

and, analogously, $\mathbb{E}[E_t] = c_t^\top \phi$. It follows that $\mathbb{E}[Y_t^\top] = x_t^\top \phi$, for any $t \in \mathbb{N}$. Now, for any $t \in \mathbb{N}$, using the very same arguments as in the proof of Theorem 1, from the fact that $\mathbb{E}[Y_t^\top] = x_t^\top \phi$ we can deduce an analogous of (1), and, from the fact that the variances of the random variables $D_1, E_1, \dots, D_t, E_t$ (for the indexes for which they are defined) are less than or equal to 1, we can deduce an analogous of (2). These two results team up to yield a bound analogous to (3): for $t \in \{2, 3, \dots\}$,

$$\mathbb{E}[|c_t^\top \hat{\phi}_{t-1} - c_t^\top \phi|^2] \leq 2 \|c_t\|_{(x_{t-1} x_{t-1}^\top + d^{-1} \mathbf{1}_d)^{-1}}^2.$$

Hence, leveraging Corollary 1, for any $T \in \mathbb{N}$, we have that

$$\begin{aligned} R_T &\leq \sum_{t=1}^T 1 \wedge \left(L \mathbb{E}[|P_t - c_t^\top \phi|^2] \right) \leq \sum_{t=2}^T (1 - b_t) L \mathbb{E}[|c_t^\top \hat{\phi}_{t-1} - c_t^\top \phi|^2] + \sum_{t=1}^T b_t \\ &\leq L \sum_{t=2}^T (1 - b_t) \|\sqrt{2}c_t\|_{(x_{t-1} x_{t-1}^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 + \sum_{t=1}^T b_t \leq \sqrt{2LdT \ln(1 + 2d(T - 1))} + \sum_{t=1}^T b_t. \end{aligned}$$

Now, given that $LT/(2d\ln(1 + 2d(T - 1))) \geq 1$, using the convention $0/0 = 0$,

$$\begin{aligned} \sum_{t=2}^T b_t &= \sum_{t=2}^T \frac{b_t \|\sqrt{2}c_t\|_{(x_{t-1} x_{t-1}^\top + d^{-1} \mathbf{1}_d)^{-1}}^2}{\|\sqrt{2}c_t\|_{(x_{t-1} x_{t-1}^\top + d^{-1} \mathbf{1}_d)^{-1}}^2} \leq \sqrt{\frac{LT}{2d\ln(1 + 2d(T - 1))}} \sum_{t=2}^T 1 \wedge b_t \|\sqrt{2}c_t\|_{(2 \sum_{s=1}^{t-1} b_s c_s c_s^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 \\ &= \sqrt{LT/(2d\ln(1 + 2d(T - 1)))} \sum_{t=1}^{T-1} 1 \wedge \|b_{t+1} \sqrt{2}c_{t+1}\|_{(\sum_{s=1}^t (b_s \sqrt{2}c_s)(b_s \sqrt{2}c_s)^\top + d^{-1} \mathbf{1}_d)^{-1}}^2 = (*). \end{aligned}$$

Using the elliptical potential lemma [37, Lemma 19.4], we obtain

$$\sum_{t=1}^T b_t \leq 1 + (*) \leq 1 + \sqrt{LT/(2d\ln(1 + 2d(T - 1)))} \cdot 2d\ln(1 + 2d(T - 1)) = 1 + \sqrt{2LdT \ln(1 + 2d(T - 1))}.$$

Hence, if $d < T/2$, this implies that $R_T \leq 1 + 2\sqrt{2LdT \ln(1 + 2d(T - 1))} \leq 1 + 4\sqrt{LdT \ln T}$. On the other hand, if $d \geq T/2$, then, since $L \geq 1$, we obtain, again, $R_T \leq T \leq 1 + 4\sqrt{LdT \ln T}$. \square

We conclude this section by stating a matching (up to logarithmic terms) worst-case $\Omega(\sqrt{LdT})$ regret lower bound for any algorithm in the two-bit-feedback case, proving the optimality of Algorithm 2.

At a high level, the proof of this result is based on the same trick (as in the proof of Theorem 2) of choosing contexts equal to vectors of the canonical basis of \mathbb{R}^d in order to obtain d independent 1-dimensional sub-instances. In each block, an adaptation of the non-contextual full-feedback lower bound construction in [9, Theorem 5] yields a lower bound of order $\sqrt{LT/d}$. Summing over blocks gives the result. For more details on the proof of this result, see Appendix B.

Theorem 4. There exist two numerical constants $a, b > 0$ such that, for any $L \geq 2$ and any time horizon $T \geq \max(4, adL^3, 2d)$, there exists a sequence of contexts $c_1, \dots, c_T \in [0, 1]^d$ such that, for any algorithm α for the contextual brokerage problem with two-bit feedback, there exists a vector $\phi \in [0, 1]^d$ and two zero-mean independent sequences $(\xi_t)_{t \in [T]}$ and $(\zeta_t)_{t \in [T]}$ independent of each other such that, if we define $V_t := c_t^\top \phi + \xi_t$ and $W_t := c_t^\top \phi + \zeta_t$, then for each $t \in [T]$ it holds that $c_t^\top \phi \in [0, 1]$, V_t and W_t are $[0, 1]$ -valued random variables with density bounded by L , and the regret of α on the sequence of traders' valuations $V_1, W_1, \dots, V_T, W_T$ satisfies $R_T \geq b\sqrt{LdT}$.

We remark that the previous lower bound holds even for algorithms that have prior knowledge of the sequence of contexts c_1, c_2, \dots and that Theorem 3 shows that Algorithm 2 matches the optimal \sqrt{LdT} rate (up to a $\sqrt{\ln T}$ factor) even without this *a-priori* knowledge.

5 Beyond Bounded Densities

In this final section, we investigate the general case where the valuations of the traders are not assumed to have a bounded density, and we show that the problem is, in general, unlearnable.

At a high level, the main reason why the problem becomes unlearnable is that Lemma 1 and its Corollary 1 fail to hold. In fact, the optimal price at time t depends in general not only on the market value $m_t = c_t^\top \phi$, but also on properties of the *time-varying* distributions of the perturbations ξ_t and ζ_t , which essentially turns our problem into a fully-adversarial one where we strive to compete against time-varying policies. For a full proof of the following theorem, see Appendix C.

Theorem 5. *There exists a sequence of contexts $c_1, c_2, \dots \in [0, 1]^d$ and a vector $\phi \in [0, 1]^d$, such that for any algorithm α for the contextual brokerage problem under full feedback, there exists an independent sequence of zero mean random variables $\xi_1, \zeta_1, \xi_2, \zeta_2, \dots$, such that if the valuations of the traders at time t are $V_t = c_t^\top \phi + \xi_t$ and $W_t = c_t^\top \phi + \zeta_t$, then $c_t^\top \phi \in [0, 1]$, V_t, W_t are $[0, 1]$ -valued random variables, and the regret of α on the sequence of traders' valuations $V_1, W_1, \dots, V_T, W_T$ satisfies $R_T = \Omega(T)$.*

We remark that the previous unlearnability result holds even for algorithms that have prior knowledge of the sequence of contexts c_1, c_2, \dots and, strikingly, of the vector ϕ .

6 Conclusions, Limitations, and Open Problems

Conclusions. Motivated by the real-life *desideratum* to exploit prior information on the traded assets, we investigated the noisy linear contextual online learning problem of brokerage between traders without predetermined seller/buyer roles. We provided a complete picture with tight regret bounds in all the proposed settings, i.e., under full and two-bit feedback, and with or without regularity assumptions on the noise distributions, achieving tightness (up to log terms) in all relevant parameters.

Limitations and Open Problems. One limitation of our work is the assumption that contexts have a linear relationship with market values. While linearity is a common and reasonable assumption in contextual settings, exploring other relationships beyond the linear case would be interesting. Another limitation of our setting is that it models traders' valuations as noisy perturbations of the market value and does not encompass cases where contextual information on traders' profiles is available. We leave the study of these extensions to future research.

Acknowledgements

The work of FB was supported by the Projects BOLD (ANR-19-CE23-0026) and GAP (ANR-21-CE40-0007) of the French National Research Agency (ANR). RC is partially supported by the MUR PRIN grant 2022EKNE5K (Learning in Markets and Society), funded by the NextGenerationEU program within the PNRR scheme, the FAIR (Future Artificial Intelligence Research) project, funded by the NextGenerationEU program within the PNRR-PE-AI scheme, the EU Horizon CL4-2022-HUMAN-02 research and innovation action under grant agreement 101120237, project ELIAS (European Lighthouse of AI for Sustainability). 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 DGECR2023-00208 (Discovery Grants Program, DGECR - Discovery Launch Supplement).

References

- [1] OTC derivatives statistics at end-June 2022. *Bank for International Settlements*, 2022.
- [2] T. Archbold, B. de Keijzer, and C. Ventre. Non-obvious manipulability for single-parameter agents and bilateral trade. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pages 2107–2115, USA, 2023. International Foundation for Autonomous Agents and Multiagent Systems.
- [3] Y. Azar, A. Fiat, and F. Fusco. An alpha-regret analysis of adversarial bilateral trade. *Advances in Neural Information Processing Systems*, 35:1685–1697, 2022.
- [4] M. Babaioff, K. Goldner, and Y. 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, page 2452–2471, USA, 2020. Society for Industrial and Applied Mathematics.
- [5] A. Badanidiyuru, Z. Feng, and G. Guruganesh. Learning to bid in contextual first price auctions. In *Proceedings of the ACM Web Conference 2023*, pages 3489–3497, 2023.
- [6] R. F. Bass. *Real analysis for graduate students*. Createspace Ind Pub, USA, 2013.
- [7] M. Bernasconi, M. Castiglioni, A. Celli, and F. Fusco. No-regret learning in bilateral trade via global budget balance. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, 2024.
- [8] L. Blumrosen and Y. Mizrahi. Approximating gains-from-trade in bilateral trading. In *Web and Internet Economics, WINE’16*, volume 10123 of *Lecture Notes in Computer Science*, pages 400–413, Germany, 2016. Springer.
- [9] N. Bolić, T. Cesari, and R. Colomboni. An online learning theory of brokerage. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, AAMAS ’24, page 216–224, Richland, SC, 2024. International Foundation for Autonomous Agents and Multiagent Systems.
- [10] J. Brustle, Y. Cai, F. Wu, and M. Zhao. Approximating gains from trade in two-sided markets via simple mechanisms. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, EC ’17, page 589–590, New York, NY, USA, 2017. Association for Computing Machinery.
- [11] J. Bu, D. Simchi-Levi, and C. Wang. Context-based dynamic pricing with partially linear demand model. *Advances in Neural Information Processing Systems*, 35:23780–23791, 2022.
- [12] N. Cesa-Bianchi, T. Cesari, R. Colomboni, F. Fusco, and S. Leonardi. Bilateral trade: A regret minimization perspective. *Mathematics of Operations Research*, 49(1):171–203, 2024.
- [13] N. Cesa-Bianchi, T. Cesari, R. Colomboni, F. Fusco, and S. Leonardi. Regret analysis of bilateral trade with a smoothed adversary. Technical report hal-04383576, 2024.
- [14] N. Cesa-Bianchi, T. R. Cesari, R. Colomboni, F. Fusco, and S. Leonardi. A regret analysis of bilateral trade. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pages 289–309, USA, 2021. Association for Computing Machinery.
- [15] N. Cesa-Bianchi, T. R. Cesari, R. Colomboni, F. Fusco, and S. Leonardi. Repeated bilateral trade against a smoothed adversary. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 1095–1130, USA, 2023. PMLR, PMLR.
- [16] T. R. Cesari and R. Colomboni. A nearest neighbor characterization of Lebesgue points in metric measure spaces. *Mathematical Statistics and Learning*, 3(1):71–112, 2021.
- [17] N. Chen and G. Gallego. Nonparametric pricing analytics with customer covariates. *Operations Research*, 69(3):974–984, 2021.

- [18] M. C. Cohen, I. Lobel, and R. Paes Leme. Feature-based dynamic pricing. In *Proceedings of the 2016 ACM Conference on Economics and Computation*, pages 817–817, 2016.
- [19] M. C. Cohen, I. Lobel, and R. Paes Leme. Feature-based dynamic pricing. *Management Science*, 66(11):4921–4943, 2020.
- [20] R. Colini-Baldeschi, B. de Keijzer, S. Leonardi, and S. Turchetta. Approximately efficient double auctions with strong budget balance. In *ACM-SIAM Symposium on Discrete Algorithms, SODA’16*, pages 1424–1443, USA, 2016. SIAM.
- [21] R. Colini-Baldeschi, P. W. Goldberg, B. de Keijzer, S. Leonardi, and S. Turchetta. Fixed price approximability of the optimal gain from trade. In *Web and Internet Economics, WINE’17*, volume 10660 of *Lecture Notes in Computer Science*, pages 146–160, Germany, 2017. Springer.
- [22] R. Colini-Baldeschi, P. W. Goldberg, B. d. Keijzer, S. Leonardi, T. Roughgarden, and S. Turchetta. Approximately efficient two-sided combinatorial auctions. *ACM Transactions on Economics and Computation (TEAC)*, 8(1):1–29, 2020.
- [23] Y. Deng, J. Mao, B. Sivan, and K. Wang. Approximately efficient bilateral trade. In *STOC*, pages 718–721, Italy, 2022. ACM.
- [24] P. Dütting, F. Fusco, P. Lazos, S. Leonardi, and R. Reiffenhäuser. Efficient two-sided markets with limited information. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2021*, page 1452–1465, New York, NY, USA, 2021. Association for Computing Machinery.
- [25] J. Fan, Y. Guo, and M. Yu. Policy optimization using semiparametric models for dynamic pricing. *Journal of the American Statistical Association*, 119(545):552–564, 2024.
- [26] D. Foster, A. Rakhlin, D. Simchi-Levi, and Y. Xu. Instance-dependent complexity of contextual bandits and reinforcement learning: A disagreement-based perspective. In *Conference on Learning Theory*, pages 2059–2059. PMLR, 2021.
- [27] D. J. Foster and A. Krishnamurthy. Contextual bandits with surrogate losses: Margin bounds and efficient algorithms. *Advances in Neural Information Processing Systems*, 31, 2018.
- [28] D. J. Foster, A. Krishnamurthy, and H. Luo. Model selection for contextual bandits. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [29] O. A. Hanna, L. Yang, and C. Fragouli. Contexts can be cheap: Solving stochastic contextual bandits with linear bandit algorithms. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 1791–1821. PMLR, 2023.
- [30] Y. Hu, N. Kallus, and X. Mao. Smooth contextual bandits: Bridging the parametric and non-differentiable regret regimes. In *Conference on Learning Theory*, pages 2007–2010. PMLR, 2020.
- [31] A. Javanmard and H. Nazerzadeh. Dynamic pricing in high-dimensions. *Journal of Machine Learning Research*, 20(9):1–49, 2019.
- [32] Z. Y. Kang, F. Pernice, and J. Vondrák. Fixed-price approximations in bilateral trade. In *Proceedings of the 2022 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 2964–2985, Alexandria, VA, USA, 2022. SIAM, Society for Industrial and Applied Mathematics.
- [33] S. Kannan, J. H. Morgenstern, A. Roth, B. Waggoner, and Z. S. Wu. A smoothed analysis of the greedy algorithm for the linear contextual bandit problem. *Advances in neural information processing systems*, 31, 2018.

- [34] J. Kirschner and A. Krause. Stochastic bandits with context distributions. *Advances in Neural Information Processing Systems*, 32, 2019.
- [35] A. Krishnamurthy, J. Langford, A. Slivkins, and C. Zhang. Contextual bandits with continuous actions: Smoothing, zooming, and adapting. *Journal of Machine Learning Research*, 21(137):1–45, 2020.
- [36] A. Krishnamurthy, T. Lykouris, C. Podimata, and R. Schapire. Contextual search in the presence of irrational agents. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 910–918, 2021.
- [37] T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [38] R. P. Leme, C. Podimata, and J. Schneider. Corruption-robust contextual search through density updates. In *Conference on Learning Theory*, pages 3504–3505. PMLR, 2022.
- [39] R. P. Leme and J. Schneider. Contextual search via intrinsic volumes. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 268–282. IEEE, 2018.
- [40] R. P. Leme and J. Schneider. Contextual search via intrinsic volumes. *SIAM Journal on Computing*, 51(4):1096–1125, 2022.
- [41] A. Liu, R. P. Leme, and J. Schneider. Optimal contextual pricing and extensions. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1059–1078. SIAM, 2021.
- [42] I. Lobel, R. P. Leme, and A. Vladu. Multidimensional binary search for contextual decision-making. *Operations Research*, 66(5):1346–1361, 2018.
- [43] I. Lobel, R. Paes Leme, and A. Vladu. Multidimensional binary search for contextual decision-making. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, pages 585–585, 2017.
- [44] H. Luo, C.-Y. Wei, A. Agarwal, and J. Langford. Efficient contextual bandits in non-stationary worlds. In *Conference On Learning Theory*, pages 1739–1776. PMLR, 2018.
- [45] Y. Luo, W. W. Sun, and Y. Liu. Distribution-free contextual dynamic pricing. *Mathematics of Operations Research*, 49(1):599–618, 2024.
- [46] B. Metevier, S. Giguere, S. Brockman, A. Kobren, Y. Brun, E. Brunskill, and P. S. Thomas. Offline contextual bandits with high probability fairness guarantees. *Advances in neural information processing systems*, 32, 2019.
- [47] R. B. Myerson and M. A. Satterthwaite. Efficient mechanisms for bilateral trading. *Journal of economic theory*, 29(2):265–281, 1983.
- [48] G. Neu and J. Olkhovskaya. Efficient and robust algorithms for adversarial linear contextual bandits. In *Conference on Learning Theory*, pages 3049–3068. PMLR, 2020.
- [49] M.-h. Oh and G. Iyengar. Thompson sampling for multinomial logit contextual bandits. *Advances in Neural Information Processing Systems*, 32, 2019.
- [50] V. Shah, R. Johari, and J. Blanchet. Semi-parametric dynamic contextual pricing. *Advances in Neural Information Processing Systems*, 32, 2019.
- [51] K. Sherstyuk, K. Phankitnirundorn, and M. J. Roberts. Randomized double auctions: gains from trade, trader roles, and price discovery. *Experimental Economics*, 24(4):1–40, 2020.
- [52] A. Slivkins, K. A. Sankararaman, and D. J. Foster. Contextual bandits with packing and covering constraints: A modular Lagrangian approach via regression. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 4633–4656. PMLR, 2023.

- [53] C.-Y. Wei, H. Luo, and A. Agarwal. Taking a hint: How to leverage loss predictors in contextual bandits? In *Conference on Learning Theory*, pages 3583–3634. PMLR, 2020.
- [54] P.-O. Weill. The search theory of over-the-counter markets. *Annual Review of Economics*, 12:747–773, 2020.
- [55] D. Williams. *Probability with martingales*. Cambridge university press, UK, 1991.
- [56] J. Xu and Y.-X. Wang. Logarithmic regret in feature-based dynamic pricing. *Advances in Neural Information Processing Systems*, 34:13898–13910, 2021.
- [57] Z. Zhou, R. Xu, and J. Blanchet. Learning in generalized linear contextual bandits with stochastic delays. *Advances in Neural Information Processing Systems*, 32, 2019.

A Proof of Theorem 2

Without loss of generality, we assume that d divides T . In fact, if we prove the theorem for this case, then, by leveraging that $T \geq 2d$ and $T \geq 4$, the general case follows from

$$R_T \geq bLd \ln(\lfloor T/d \rfloor d) \geq \frac{b}{2} Ld \ln T.$$

Let $n := T/d$. Let e_1, \dots, e_d be the canonical basis of \mathbb{R}^d . Define, for all $i \in [d]$ and $j \in [n]$, the context $c_{j+(i-1)n} := e_i$. We assume that these contexts are known to the learner in advance and, therefore, we can restrict the proof to deterministic algorithms without any loss of generality.

Let $L \geq 2$, $J_L := [\frac{1}{2} - \frac{1}{14L}, \frac{1}{2} + \frac{1}{14L}]$, $f := \mathbb{I}_{[0, \frac{3}{7}]} + L\mathbb{I}_{J_L} + \mathbb{I}_{[\frac{4}{7}, 1]}$, and, for any $\varepsilon \in [-1, 1]$, $g_\varepsilon := -\varepsilon\mathbb{I}_{[\frac{1}{7}, \frac{3}{14}]} + \varepsilon\mathbb{I}_{(\frac{3}{14}, \frac{2}{7}]}$ and $f_\varepsilon := f + g_\varepsilon$. For any $\varepsilon \in [-1, 1]$, note that $0 \leq f_\varepsilon \leq L$ and $\int_0^1 f_\varepsilon(x) dx = 1$, hence f_ε is a valid density on $[0, 1]$ bounded by L . We will denote the corresponding probability measure by ν_ε , set $\bar{\nu}_\varepsilon := \int_{[0,1]} x d\nu_\varepsilon(x)$, and notice that direct computations show that $\bar{\nu}_\varepsilon = \frac{1}{2} + \frac{\varepsilon}{196}$. Consider for each $q \in [0, 1]$, an i.i.d. sequence $(B_{q,t})_{t \in \mathbb{N}}$ of Bernoulli random variables of parameter q , an i.i.d. sequence $(\tilde{B}_t)_{t \in \mathbb{N}}$ of Bernoulli random variables of parameter $1/7$, an i.i.d. sequence $(U_t)_{t \in \mathbb{N}}$ of uniform random variables on $[0, 1]$, and uniform random variables E_1, \dots, E_d on $[-\bar{\varepsilon}_L, \bar{\varepsilon}_L]$, where $\bar{\varepsilon}_L := \frac{7}{L}$, such that $((B_{q,t})_{t \in \mathbb{N}, q \in [0,1]}, (\tilde{B}_t)_{t \in \mathbb{N}}, (U_t)_{t \in \mathbb{N}}, E_1, \dots, E_d)$ is an independent family. Let $\varphi: [0, 1] \rightarrow [0, 1]$ be such that, if U is a uniform random variable on $[0, 1]$, then the distribution of $\varphi(U)$ has density $\frac{7}{6} \cdot f \cdot \mathbb{I}_{[0,1] \setminus [1/7, 2/7]}$ (which exists by the Skorokhod representation theorem [55, Section 17.3]). For each $\varepsilon \in [-1, 1]$ and $t \in \mathbb{N}$, define

$$G_{\varepsilon,t} := \left(\frac{2+U_t}{14} (1 - B_{\frac{1+\varepsilon}{2},t}) + \frac{3+U_t}{14} B_{\frac{1+\varepsilon}{2},t} \right) \tilde{B}_t + \varphi(U_t) (1 - \tilde{B}_t), \quad (4)$$

$V_{\varepsilon,t} := G_{\varepsilon,2t-1}$, $W_{\varepsilon,t} := G_{\varepsilon,2t}$, $\xi_{\varepsilon,t} := V_{\varepsilon,t} - \bar{\nu}_\varepsilon$, and $\zeta_{\varepsilon,t} := W_{\varepsilon,t} - \bar{\nu}_\varepsilon$. In the following, if a_1, \dots, a_d is a sequence of elements, we will use the notation $a_{1:d}$ as a shorthand for (a_1, \dots, a_d) . For each $\varepsilon_1, \dots, \varepsilon_d \in [-1, 1]$, each $i \in [d]$, and each $j \in [n]$, define the random variables $\xi_{j+(i-1)n}^{\varepsilon_{1:d}} := \xi_{\varepsilon_i, j+(i-1)n}$ and $\zeta_{j+(i-1)n}^{\varepsilon_{1:d}} := \zeta_{\varepsilon_i, j+(i-1)n}$. The family $(\xi_t^{\varepsilon_{1:d}}, \zeta_t^{\varepsilon_{1:d}})_{t \in [T], \varepsilon_{1:d} \in [-1,1]^d}$ is an independent family, independent of (E_1, \dots, E_d) , and for each $i \in [d]$ and each $j \in [n]$ it can be checked that the two random variables $\xi_{j+(i-1)n}^{\varepsilon_{1:d}}, \zeta_{j+(i-1)n}^{\varepsilon_{1:d}}$ are zero mean with common distribution given by ν_{ε_i} . For each $\varepsilon_1, \dots, \varepsilon_d \in [-1, 1]$, let $\phi_{\varepsilon_{1:d}} := (\bar{\nu}_{\varepsilon_1}, \dots, \bar{\nu}_{\varepsilon_d})$, and for each $i \in [d]$ and $j \in [n]$, let $V_{j+(i-1)n}^{\varepsilon_{1:d}} := c_{j+(i-1)n}^\top \phi_{\varepsilon_{1:d}} + \xi_{j+(i-1)n}^{\varepsilon_{1:d}}$ and $W_{j+(i-1)n}^{\varepsilon_{1:d}} := c_{j+(i-1)n}^\top \phi_{\varepsilon_{1:d}} + \zeta_{j+(i-1)n}^{\varepsilon_{1:d}}$. Note that these last two random variables are $[0, 1]$ -valued zero-mean perturbations of $c_{j+(i-1)n}^\top \phi_{\varepsilon_{1:d}}$ with shared density given by f_{ε_i} , and hence bounded by L .

We will show that any algorithm has to suffer the regret inequality in the statement of the theorem if the sequence of evaluations is $V_1^{\varepsilon_{1:d}}, W_1^{\varepsilon_{1:d}}, \dots, V_T^{\varepsilon_{1:d}}, W_T^{\varepsilon_{1:d}}$, for some $\varepsilon_1, \dots, \varepsilon_d \in [0, 1]$.

Before doing that, we first need the following. For any $\varepsilon_1, \dots, \varepsilon_d \in [-1, 1]$, $p \in [0, 1]$, and $t \in [T]$ let $\text{GFT}_t^{\varepsilon_{1:d}}(p) := g(p, V_t^{\varepsilon_{1:d}}, W_t^{\varepsilon_{1:d}})$.

By Lemma 1, we have, for all $\varepsilon_1, \dots, \varepsilon_d \in [-1, 1]$, $i \in [d]$, $j \in [n]$, and $p \in [0, 1]$,

$$\mathbb{E}[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(p)] = 2 \int_0^p \int_0^\lambda f_{\varepsilon_i}(s) ds d\lambda + 2(\bar{\nu}_{\varepsilon_i} - p) \int_0^p f_{\varepsilon_i}(s) ds,$$

which, together with the fundamental theorem of calculus—[6, Theorem 14.16], noting that $p \mapsto \mathbb{E}[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(p)]$ is absolutely continuous with derivative defined a.e. by $p \mapsto 2(\bar{\nu}_{\varepsilon_i} - p)f_{\varepsilon_i}(p)$ —yields, for any $p \in J_L$,

$$\mathbb{E}[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\bar{\nu}_{\varepsilon_i})] - \mathbb{E}[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(p)] = L|\bar{\nu}_{\varepsilon_i} - p|^2. \quad (5)$$

Note also that for all $\varepsilon_1, \dots, \varepsilon_d \in [-\bar{\varepsilon}_L, \bar{\varepsilon}_L]$, $t \in [T]$, and $p \in [0, 1] \setminus J_L$, a direct verification shows that

$$\mathbb{E}[\text{GFT}_t^{\varepsilon_{1:d}}(p)] \leq \mathbb{E}[\text{GFT}_t^{\varepsilon_{1:d}}(1/2)]. \quad (6)$$

Fix any arbitrary deterministic algorithm for the full feedback setting $(\alpha_t)_{t \in [T]}$, i.e., (given that the contexts c_1, \dots, c_T are here fixed and declared ahead of time to the learner), a sequence of functions $\alpha_t: ([0, 1] \times [0, 1])^{t-1} \rightarrow [0, 1]$ mapping past feedback into prices (with the convention that α_1 is just a number in $[0, 1]$). For each $t \in [T]$, define $\tilde{\alpha}_t: ([0, 1] \times [0, 1])^{t-1} \rightarrow J_L$ equal to α_t whenever α_t takes values in J_L , and equal to $1/2$ otherwise. Define $Z_1 := \frac{1+E_1}{2}, \dots, Z_d := \frac{1+E_d}{2}$.

Now, note the following

$$\begin{aligned} & \sup_{\varepsilon_{1:d} \in [-\bar{\varepsilon}_L, \bar{\varepsilon}_L]^d} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\bar{\nu}_{\varepsilon_i}) - \text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\alpha_t(V_1^{\varepsilon_{1:d}}, W_1^{\varepsilon_{1:d}}, \dots, V_{j-1+(i-1)n}^{\varepsilon_{1:d}}, W_{j-1+(i-1)n}^{\varepsilon_{1:d}}))] \\ & \stackrel{(6)}{\geq} \sup_{\varepsilon_{1:d} \in [-\bar{\varepsilon}_L, \bar{\varepsilon}_L]^d} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\bar{\nu}_{\varepsilon_i}) - \text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\tilde{\alpha}_t(V_1^{\varepsilon_{1:d}}, W_1^{\varepsilon_{1:d}}, \dots, V_{j-1+(i-1)n}^{\varepsilon_{1:d}}, W_{j-1+(i-1)n}^{\varepsilon_{1:d}}))] \\ & \stackrel{\spadesuit}{=} L \sup_{\varepsilon_{1:d} \in [-\bar{\varepsilon}_L, \bar{\varepsilon}_L]^d} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|\bar{\nu}_{\varepsilon_i} - \tilde{\alpha}_t(V_1^{\varepsilon_{1:d}}, W_1^{\varepsilon_{1:d}}, \dots, V_{j-1+(i-1)n}^{\varepsilon_{1:d}}, W_{j-1+(i-1)n}^{\varepsilon_{1:d}})|^2] \\ & \geq L \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|\bar{\nu}_{E_i} - \tilde{\alpha}_t(V_1^{E_{1:d}}, W_1^{E_{1:d}}, \dots, V_{j-1+(i-1)n}^{E_{1:d}}, W_{j-1+(i-1)n}^{E_{1:d}})|^2] \\ & \stackrel{\heartsuit}{\geq} L \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|\bar{\nu}_{E_i} - \mathbb{E}[\bar{\nu}_{E_i} | V_1^{E_{1:d}}, W_1^{E_{1:d}}, \dots, V_{j-1+(i-1)n}^{E_{1:d}}, W_{j-1+(i-1)n}^{E_{1:d}}]|^2] \\ & = \frac{L}{196} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|E_i - \mathbb{E}[E_i | V_1^{E_{1:d}}, W_1^{E_{1:d}}, \dots, V_{j-1+(i-1)n}^{E_{1:d}}, W_{j-1+(i-1)n}^{E_{1:d}}]|^2] \\ & \stackrel{\clubsuit}{\geq} \frac{L}{196} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|E_i - \mathbb{E}[E_i | B_{\frac{1+E_i}{2}, 1+2(i-1)n}, \dots, B_{\frac{1+E_i}{2}, 2(j-1)+2(i-1)n}]|^2] \\ & \stackrel{\clubsuit}{=} \frac{L}{196} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|E_i - \mathbb{E}[E_i | B_{\frac{1+E_i}{2}, 1}, \dots, B_{\frac{1+E_i}{2}, 2(j-1)}]|^2] \\ & = \frac{L}{49} \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}[|Z_i - \mathbb{E}[Z_i | B_{Z_i, 1}, \dots, B_{Z_i, 2(j-1)}]|^2] \end{aligned}$$

where \spadesuit follows from (5) and the fact that $\tilde{\alpha}_t$ takes values in J_L ; \heartsuit from the fact that the minimizer of the $L^2(\mathbb{P})$ -distance from $\bar{\nu}_{E_i}$ in $\sigma(V_1^{E_{1:d}}, W_1^{E_{1:d}}, \dots, V_{j-1+(i-1)n}^{E_{1:d}}, W_{j-1+(i-1)n}^{E_{1:d}})$ is $\mathbb{E}[\bar{\nu}_{E_i} | V_1^{E_{1:d}}, W_1^{E_{1:d}}, \dots, V_{j-1+(i-1)n}^{E_{1:d}}, W_{j-1+(i-1)n}^{E_{1:d}}]$ (see, e.g., [55, Section 9.4]); \clubsuit follows from the fact that, by Equation (4) and the independence of E_i from $((B_{q,t})_{t \in \mathbb{N}, q \in [0,1]}, (\tilde{B}_t)_{t \in \mathbb{N}}, (U_t)_{t \in \mathbb{N}})$, the conditional expectation $\mathbb{E}[E_i | V_1^{E_{1:d}}, W_1^{E_{1:d}}, \dots, V_{j-1+(i-1)n}^{E_{1:d}}, W_{j-1+(i-1)n}^{E_{1:d}}]$ is a measurable function of $B_{\frac{1+E_i}{2}, 1+2(i-1)n}, \dots, B_{\frac{1+E_i}{2}, 2(j-1)+2(i-1)n}$, together with the same observation made

in ♠ about the minimization of $L^2(\mathbb{P})$ distance; and ♦ follows from the fact that the sequence $(B_{\frac{1+E_i}{2},t})_{t \in \mathbb{N}}$ is i.i.d..

Finally, the general term of this last sum is the expected squared distance between the random parameter (drawn uniformly over $[(1 - \bar{\varepsilon}_L)/2, (1 + \bar{\varepsilon}_L)/2]$) of an i.i.d. sequence of Bernoulli random variables and the conditional expectation of this random parameter given $2(j-1)$ independent realizations of these Bernoullis. A probabilistic argument shows that there exist two universal constants $\tilde{a}, \tilde{b} > 0$ such that, for all $j \geq \tilde{b}L^4$ and each $i \in [d]$,

$$\mathbb{E}\left[\|Z_i - \mathbb{E}[Z_i | B_{Z_i,1}, \dots, B_{Z_i,2(j-1)}]\|^2\right] \geq \tilde{a} \frac{1}{j-1}. \quad (7)$$

At a high level, this is because, in an event of probability $\Omega(1)$, if j is large enough, the conditional expectation $\mathbb{E}[Z_i | B_{Z_i,1}, \dots, B_{Z_i,2(j-1)}]$ is very close to the empirical average $\frac{1}{2(j-1)} \sum_{s=1}^{2(j-1)} B_{Z_i,s}$, whose expected squared distance from Z is $\Omega(1/(j-1))$. For a formal proof of (7) with explicit constants, we refer the reader to [9, Appendix B of the extended arxiv version]. Summing over $i \in [d]$ and $j \in [n]$, we obtain that there exist $\varepsilon_1, \dots, \varepsilon_d \in [-1, 1]^d$ such that

$$\begin{aligned} & \sum_{i=1}^d \sum_{j=1}^n \mathbb{E}\left[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\bar{\nu}_{\varepsilon_i}) - \text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(\tilde{\alpha}_t(V_1^{\varepsilon_{1:d}}, W_1^{\varepsilon_{1:d}}, \dots, V_{j-1+(i-1)n}^{\varepsilon_{1:d}}, W_{j-1+(i-1)n}^{\varepsilon_{1:d}}))\right] \\ &= \Omega(Ld \ln n) = \Omega(Ld \ln T). \end{aligned}$$

B Proof of Theorem 4

Fix $L \geq 2$ and $T \in \mathbb{N}$. We will use the very same notation as in the proof of Theorem 2. In particular, the contexts c_1, \dots, c_T are again the same as before and declared ahead of time to the learner. We will show that for each algorithm for contextual brokerage with 2-bit feedback and each time horizon T , if $R_T^{\varepsilon_{1:d}}$ is the regret of the algorithm at time horizon T when the traders' valuations are $V_1^{\varepsilon_{1:d}}, W_1^{\varepsilon_{1:d}}, \dots, V_T^{\varepsilon_{1:d}}, W_T^{\varepsilon_{1:d}}$, then $\max_{\sigma_{1:d} \in \{-1, 1\}^d} R_T^{(\sigma_1 \varepsilon, \dots, \sigma_d \varepsilon)} = \Omega(\sqrt{dT})$ if $\varepsilon = \Theta((LT/d)^{-1/4})$ and $T = \Omega(dL^3)$.

Note that for all $\varepsilon_{1:d} \in [-1, 1]^d$, $i \in [d]$, $j \in [n]$, and $p < \frac{1}{2}$, if $\varepsilon_i > 0$, then, a direct verification shows that

$$\mathbb{E}\left[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(1/2)\right] \geq \mathbb{E}\left[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(p)\right]. \quad (8)$$

Similarly, for all $\varepsilon_{1:d} \in [-1, 1]^d$, $i \in [d]$, $j \in [n]$, and $p > \frac{1}{2}$, if $\varepsilon_i < 0$, then

$$\mathbb{E}\left[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(1/2)\right] \geq \mathbb{E}\left[\text{GFT}_{j+(i-1)n}^{\varepsilon_{1:d}}(p)\right]. \quad (9)$$

Furthermore, a direct verification shows that, for each $\varepsilon_{1:d} \in [-1, 1]^d$ and $t \in [T]$,

$$\max_{p \in [0, 1]} \mathbb{E}\left[\text{GFT}_t^{\varepsilon_{1:d}}(p)\right] - \max_{p \in [\frac{1}{7}, \frac{2}{7}]} \mathbb{E}\left[\text{GFT}_t^{\varepsilon_{1:d}}(p)\right] \geq \frac{1}{50} = \Omega(1). \quad (10)$$

Now, assume that $T \geq dL^3/14^4$ so that, defining $\varepsilon := (LT/d)^{-1/4}$, we have that for any $\sigma_{1:d} \in \{-1, 1\}^d$, any $i \in [d]$ and any $j \in [n]$, the maximizer of the expected gain from trade $p \mapsto \mathbb{E}[\text{GFT}_{j+(i-1)n}^{(\sigma_1 \varepsilon, \dots, \sigma_d \varepsilon)}(p)]$ is at $\frac{1}{2} + \frac{\sigma_i \varepsilon}{196}$ and hence belongs to the spike region J_L . If $\sigma_i = 1$ (resp., $\sigma_i = -1$) case, the optimal price for the rounds $1 + (i-1)n, \dots, in$ belongs to the region $(\frac{1}{2}, \frac{1}{2} + \frac{1}{14L}]$ (resp., $[\frac{1}{2} - \frac{1}{14L}, \frac{1}{2})$). By posting prices in the wrong region $[0, \frac{1}{2}]$ (resp., $[\frac{1}{2}, 1]$) in the $\sigma_i = 1$ (resp., $\sigma_i = -1$) case, the learner incurs a $\Omega(L\varepsilon^2) = \Omega(\sqrt{L/dT})$ instantaneous regret by (5) and (8) (resp., (5) and (9)). Then, in order to attempt suffering less than $\Omega(\sqrt{L/T} \cdot n) = \Omega(\sqrt{LT/d})$ regret in the rounds $1 + (i-1)n, \dots, in$, the algorithm would have to detect the sign of σ_i and play accordingly. We will show now that even this strategy will not improve the regret of the algorithm (by more than a constant) because of the cost of determining the sign of σ_i with the available feedback. Since for any $i \in [d]$ and $j \in [n]$, the feedback received from the two traders at time

$j + (i - 1)n$ by posting a price p is $\mathbb{I}\{p \leq V_{j+(i-1)n}^{(\sigma_1\varepsilon, \dots, \sigma_d\varepsilon)}\}$ and $\mathbb{I}\{p \leq W_{j+(i-1)n}^{(\sigma_1\varepsilon, \dots, \sigma_d\varepsilon)}\}$, the only way to obtain information about (the sign of) σ_i is to post in the costly $(\Omega(1)$ -instantaneous regret by Equation (10)) sub-optimal region $[\frac{1}{7}, \frac{2}{7}]$. However, posting prices in the region $[\frac{1}{7}, \frac{2}{7}]$ at time $j + (i - 1)n$ can't give more information about σ_i than the information carried by $V_{j+(i-1)n}^{(\sigma_1\varepsilon, \dots, \sigma_d\varepsilon)}$ and $W_{j+(i-1)n}^{(\sigma_1\varepsilon, \dots, \sigma_d\varepsilon)}$, which, in turn, can't give more information about σ_i than the information carried by the two Bernoullis $B_{\frac{1+\sigma_i\varepsilon}{2}, 2(j+(i-1)n)-1}$ and $B_{\frac{1+\sigma_i\varepsilon}{2}, 2(j+(i-1)n)}$. Since only during rounds $1 + (i - 1)n, \dots, in$ is possible to extract information about the sign of σ_i and, (via an information-theoretic argument) in order to distinguish the sign of σ_i having access to i.i.d. Bernoulli random variables of parameter $\frac{1+\sigma_i\varepsilon}{2}$ requires $\Omega(1/\varepsilon^2) = \Omega(\sqrt{LT/d})$ samples, we are forced to post at least $\Omega(\sqrt{LT/d})$ prices in the costly region $[\frac{1}{7}, \frac{2}{7}]$ during the rounds $1 + (i - 1)n, \dots, in$ suffering a regret of $\Omega(\sqrt{LT/d}) \cdot \Omega(1) = \Omega(\sqrt{LT/d})$. Putting everything together, no matter what the strategy, each algorithm will pay at least $\Omega(\sqrt{LT/d})$ regret in each epoch $1 + (i - 1)n, \dots, in$ for every $i \in [d]$, resulting in an overall regret of $\Omega(\sqrt{LT/d}) \cdot d = \Omega(\sqrt{dLT})$.

C Proof of Theorem 5

Assume that $d \geq 2$ (for the case $d = 1$, the following proof can be adapted straightforwardly by defining $\phi = 1$ and $c_t = 1/2 + \varepsilon_t$, where ε_t is an arbitrary small sequence of biases). Let $(a_t)_{t \in \mathbb{N}}$ be a sequence of distinct elements in $[0, 1]$ and, for all $t \in \mathbb{N}$, let $c_t := (a_t, 1 - a_t, 0, 0, \dots, 0)$. Notice that $(c_t)_{t \in \mathbb{N}}$ is a sequence of distinct elements in $[0, 1]^2$. Define $\phi := (1/2, 1/2, 0, 0, \dots, 0)$. Notice that for each $t \in \mathbb{N}$ it holds that $c_t^\top \phi = 1/2$. Let $\varepsilon \in (0, 1/16)$. For any $\theta \in \{0, 1\}$, consider the following probability distribution

$$\mu_\theta := \left(\frac{1}{4} + (1 - 2\theta)\varepsilon \right) \delta_{-\frac{1}{2}} + \frac{1}{2} \delta_{2(1-\theta)\varepsilon - 2\theta\varepsilon} + \left(\frac{1}{4} - (1 - 2\theta)\varepsilon \right) \delta_{\frac{1}{2}},$$

where for any $a \in \mathbb{R}$, δ_a is the Dirac's delta probability distribution centered in a . Consider an independent family of random variables $(\xi_{t,\theta}, \zeta_{t,\theta})_{t \in \mathbb{N}, \theta \in \{0,1\}}$ such that for any $t \in \mathbb{N}$ and any $\theta \in \{0, 1\}$, we have that both $\xi_{t,\theta}$ and $\zeta_{t,\theta}$ are random variables with common distribution μ_θ . Notice that for each $t \in \mathbb{N}$ and each $\theta \in \{0, 1\}$ we have that $\mathbb{E}[\xi_{t,\theta}] = 0 = \mathbb{E}[\zeta_{t,\theta}]$. Define, for each $t \in \mathbb{N}$ and each $\theta \in \{0, 1\}$, the random variables $V_{t,\theta} := c_t^\top \phi + \xi_t$ and $W_{t,\theta} := c_t^\top \phi + \zeta_t$. Notice that these are $[0, 1]$ -valued random variables and that $(V_{t,\theta}, W_{t,\theta})_{t \in \mathbb{N}, \theta \in \{0,1\}}$ is an independent family. Now, for each $\theta \in \{0, 1\}$ and each $t \in \mathbb{N}$, let

$$p^\#(\theta) \in \operatorname{argmax}_{p \in [0,1]} \mathbb{E}[g(p, V_{t,\theta}, W_{t,\theta})],$$

which does exist because the function $[0, 1] \rightarrow [0, 1], p \mapsto \mathbb{E}[g(p, V_{t,\theta}, W_{t,\theta})]$ is upper semicontinuous (this can be proved as in [12, Appendix B]) and defined on a compact set. Furthermore, note that the previous definition is independent of t because, for any $\theta \in \{0, 1\}$, the pairs $(V_{t_1,\theta}, W_{t_1,\theta})$ and $(V_{t_2,\theta}, W_{t_2,\theta})$ share the same distribution for every $t_1, t_2 \in \mathbb{N}$. Fix a learning algorithm for the full-feedback contextual brokerage problem, fix a time horizon $T \in \mathbb{N}$, and notice that since the contexts c_1, c_2, \dots are all distinct, it follows that

$$\begin{aligned} & \max_{\theta_1, \dots, \theta_T \in \{0,1\}^T} \sup_{p^*:[0,1]^d \rightarrow [0,1]} \mathbb{E} \left[\sum_{t=1}^T (g(p^*(c_t), V_{t,\theta_t}, W_{t,\theta_t}) - g(P_t, V_{t,\theta_t}, W_{t,\theta_t})) \right] \\ &= \max_{\theta_1, \dots, \theta_T \in \{0,1\}^T} \sum_{t=1}^T \left(\sup_{p \in [0,1]} \mathbb{E}[g(p, V_{t,\theta_t}, W_{t,\theta_t})] - \mathbb{E}[g(P_t, V_{t,\theta_t}, W_{t,\theta_t})] \right) \\ &= \max_{\theta_1, \dots, \theta_T \in \{0,1\}^T} \mathbb{E} \left[\sum_{t=1}^T (g(p^\#(\theta_t), V_{t,\theta_t}, W_{t,\theta_t}) - g(P_t, V_{t,\theta_t}, W_{t,\theta_t})) \right] = (\#). \end{aligned}$$

Now, consider an i.i.d. family of Bernoulli random variables $(\Theta_t)_{t \in \mathbb{N}}$ with parameter 1/2, independent of the whole family $(V_{t,\theta}, W_{t,\theta})_{t \in \mathbb{N}, \theta \in \{0,1\}}$. We have that

$$\begin{aligned} (\#) &\geq \mathbb{E} \left[\sum_{t=1}^T \left(g(p^\#(\Theta_t), V_{t,\Theta_t}, W_{t,\Theta_t}) - g(P_t, V_{t,\Theta_t}, W_{t,\Theta_t}) \right) \right] \\ &= \sum_{t=1}^T \left(\mathbb{E} \left[g(p^\#(\Theta_t), V_{t,\Theta_t}, W_{t,\Theta_t}) \right] - \mathbb{E} \left[g(P_t, V_{t,\Theta_t}, W_{t,\Theta_t}) \right] \right) =: (\$) \end{aligned}$$

Now, for each $t \in [T]$, we see that

$$\begin{aligned} \mathbb{E} \left[g(p^\#(\Theta_t), V_{t,\Theta_t}, W_{t,\Theta_t}) \right] &= \mathbb{E} \left[\mathbb{E} \left[g(p^\#(\Theta_t), V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t \right] \right] \\ &= \mathbb{E} \left[\max_{p \in [0,1]} \mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t \right] \right] \end{aligned}$$

and long but straightforward computations show that, for each $p \in [0,1]$, it holds that

$$\mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t \right] = \begin{cases} \frac{1}{4} + \varepsilon(1 - 2\Theta_t) & \text{if } 0 \leq p < \frac{1}{2} - 2\Theta_t\varepsilon + 2(1 - \Theta_t)\varepsilon, \\ \frac{3}{8} + 2\varepsilon^2 & \text{if } p = \frac{1}{2} - 2\Theta_t\varepsilon + 2(1 - \Theta_t)\varepsilon, \\ \frac{1}{4} - \varepsilon(1 - 2\Theta_t) & \text{if } \frac{1}{2} - 2\Theta_t\varepsilon + 2(1 - \Theta_t)\varepsilon < p \leq 1, \end{cases}$$

from which it follows that

$$\max_{p \in [0,1]} \mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t \right] = \frac{3}{8} + 2\varepsilon^2.$$

On the other hand, for each $t \in [T]$, leveraging the freezing lemma [16, Lemma 8], we have that

$$\begin{aligned} \mathbb{E} \left[g(P_t, V_{t,\Theta_t}, W_{t,\Theta_t}) \right] &= \mathbb{E} \left[\mathbb{E} \left[g(P_t, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid P_t \right] \right] = \mathbb{E} \left[\left[\mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \right] \right]_{p=P_t} \right] \\ &= \mathbb{E} \left[\left[\frac{1}{2} \mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t = 0 \right] + \frac{1}{2} \mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t = 1 \right] \right]_{p=P_t} \right] \end{aligned}$$

and again, tedious but straightforward computations show that, for each $p \in [0,1]$, it holds that

$$\begin{aligned} &\frac{1}{2} \mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t = 0 \right] + \frac{1}{2} \mathbb{E} \left[g(p, V_{t,\Theta_t}, W_{t,\Theta_t}) \mid \Theta_t = 1 \right] \\ &= \frac{1}{4} \left(\mathbb{I} \left\{ p < \frac{1}{2} - 2\varepsilon \right\} + \mathbb{I} \left\{ \frac{1}{2} + 2\varepsilon < p \right\} \right) + \left(\frac{5}{16} + \frac{\varepsilon}{2} + \varepsilon^2 \right) \left(\mathbb{I} \left\{ p = \frac{1}{2} - 2\varepsilon \right\} + \mathbb{I} \left\{ p = \frac{1}{2} + 2\varepsilon \right\} \right) \\ &\quad + \left(\frac{1}{4} + \varepsilon \right) \mathbb{I} \left\{ \frac{1}{2} - 2\varepsilon < p < \frac{1}{2} + 2\varepsilon \right\} \\ &\leq \frac{5}{16} + \frac{\varepsilon}{2} + \varepsilon^2. \end{aligned}$$

We conclude that

$$(\$) \geq \frac{T}{16} + \left(\varepsilon^2 - \frac{\varepsilon}{2} \right) T,$$

from which it follows that there exists $\theta_1, \dots, \theta_T \in \{0,1\}$ such that

$$\sup_{p^* : [0,1]^d \rightarrow [0,1]} \mathbb{E} \left[\sum_{t=1}^T \left(g(p^*(c_t), V_{t,\theta_t}, W_{t,\theta_t}) - g(P_t, V_{t,\theta_t}, W_{t,\theta_t}) \right) \right] \geq \frac{T}{16} + \left(\varepsilon^2 - \frac{\varepsilon}{2} \right) T \geq \frac{T}{32}.$$