# No-Regret Learning in Bilateral Trade via Global Budget Balance[*]

Martino Bernasconi[†]    Matteo Castiglioni[‡]    Andrea Celli[†]    Federico Fusco[#]

[†] Bocconi university
[‡] Politecnico di Milano
[#] Sapienza University of Rome

{martino.bernasconi,andrea.celli2}@unibocconi.it,  matteo.castiglioni@polimi.it,
federico.fusco@uniroma1.it

## Abstract

Bilateral trade models the problem of intermediating between two rational agents — a seller and a buyer — both characterized by a private valuation for an item they want to trade. We study the online learning version of the problem, in which at each time step a new seller and buyer arrive and the learner has to set prices for them without any knowledge about their (adversarially generated) valuations.

In this setting, known impossibility results rule out the existence of no-regret algorithms when budget balanced has to be enforced at each time step. In this paper, we introduce the notion of *global budget balance*, which only requires the learner to fulfill budget balance over the entire time horizon. Under this natural relaxation, we provide the first no-regret algorithms for adversarial bilateral trade under various feedback models. First, we show that in the full-feedback model, the learner can guarantee $\tilde{O}(\sqrt{T})$ regret against the best fixed prices in hindsight, and that this bound is optimal up to poly-logarithmic terms. Second, we provide a learning algorithm guaranteeing a $\tilde{O}(T^{3/4})$ regret upper bound with one-bit feedback, which we complement with a $\Omega(T^{5/7})$ lower bound that holds even in the two-bit feedback model. Finally, we introduce and analyze an alternative benchmark that is provably stronger than the best fixed prices in hindsight and is inspired by the literature on bandits with knapsacks.

arXiv:2310.12370v2 [cs.GT] 27 Mar 2024

# Contents

# 1 Introduction

Bilateral trade is a classic economic problem where two agents — a seller and a buyer — are interested in trading a good. Both agents are characterized by a private valuation for the item, and their goal is to maximize their own utility. Solving this problem requires the design of a mechanism that intermediates between the two parties, facilitating the trade. Ideally, the mechanism should maximize efficiency (*i.e.,* trade whenever the buyer's valuation exceeds the seller's one) while ensuring that agents behave according to their true preferences (*incentive compatibility*), and that the utility for participating in the mechanism of each agent is non-negative (*individual rationality*). These properties ensure favorable outcomes for the agents, yet they do not guarantee the economic viability of the mechanism. To see this, consider the following mechanism $\mathcal{M}$. $\mathcal{M}$ asks the agents for their valuations, $s$ for the seller and $b$ for the buyer, and makes the trade happen if it is convenient (*i.e.,* if $s \leq b$). In case of a trade, $\mathcal{M}$ then charges $s$ to the buyer and pays $b$ to the buyer. It is not hard to see that $\mathcal{M}$ enforces incentive compatibility and individual rationality, and is efficient by design. However, it exhibits the major drawback of allowing the intermediary to incur a net loss when $b > s$. To avoid such situations, a crucial constraint in bilateral trade is *budget balance*, which restricts the mechanism from subsidizing the agents.

As highlighted by the above example, an incentive compatible mechanism maximizing efficiency for bilateral trade may not be budget balanced. This phenomenon was first observed by Vickrey [1961]; subsequently Myerson and Satterthwaite [1983], provided a more general impossibility result by showing the existence of instances where a fully efficient mechanism that satisfies incentive compatibility, individual rationality, and budget balance does not exist. This result holds even when probabilistic information on the agents' valuations is available. To circumvent these impossibility results, the extensive subsequent research primarily focuses on finding approximately efficient mechanisms in the Bayesian setting. There, various incentive compatible mechanisms exist that give a constant-factor approximation to the social welfare (see, *e.g.,* Blumrosen and Dobzinski [2014], Kang et al. [2022], while more recent works also consider the harder problem of approximating the gain from trade [McAfee, 2008, Blumrosen and Mizrahi, 2016, Brustle et al., 2017, Deng et al., 2022, Fei, 2022]. While the Bayesian assumption of having perfect knowledge about the underlying distributions of valuations is, in some sense, necessary for extracting meaningful approximations to the social welfare [Dütting et al., 2021], it is important to observe that this assumption is oftentimes unrealistic.

Following the recent line of work initiated by Cesa-Bianchi et al. [2021], we study this fundamental mechanism design problem through the lens of regret minimization in a repeated setting where at each time $t$, a new seller/buyer pair arrives. The seller arriving at time $t$ has a private valuation $s_t$ representing the lowest price they are willing to accept for the item. Analogously, the buyer has a private valuation $b_t$ representing the highest price they are willing to pay for the item. The learner, without any knowledge about the private valuations at the current time $t$, posts two (possibly randomized) prices: $p_t$ to the seller and $q_t$ to the buyer. A trade happens when both agents agree to trade, i.e., when $s_t \leq p_t$ and $q_t \leq b_t$. After posting $(p_t, q_t)$, the learner observes some feedback about the transaction, and is awarded the *gain from trade*:

$$\text{GFT}_t(p, q) = \mathbb{I}\{s_t \leq p\}\mathbb{I}\{q \leq b_t\}(b_t - s_t).$$

The goal of the learner is to maximize the overall gain from trade or, equivalently, minimize the regret with respect to the best price in hindsight. Prior research has investigated the impact of different budget balance notions on the problem's learnability. When the mechanism is constrained to enforce per-round *strong budget balance* (*i.e.,* $p_t = q_t$ at each time step $t$), it is possible to attain sublinear regret only when the sequence of valuations is drawn i.i.d. from some fixed unknown

| | Type of Adversary | Budget Balance | Regret Upper Bounds | Regret Lower Bounds |
|---|---|---|---|---|
| Cesa-Bianchi et al. [2021] | stochastic setting | strong | • Full: $\tilde{O}(T^{1/2})$ <br> • Partial: $\tilde{O}(T^{2/3})^*$ | • Full: $\Omega(T^{1/2})$ <br> • Partial: $\Omega(T^{2/3})$ |
| Azar et al. [2022] | adversarial setting | weak | • Full: $\tilde{O}(T^{1/2})^\dagger$ <br> • Partial: $\tilde{O}(T^{3/4})^\dagger$ | • Full: $\Omega(T^{1/2})^\dagger$ <br> • Partial: $\Omega(T^{2/3})^\dagger$ |
| Cesa-Bianchi et al. [2023] | $\sigma$-smooth adversary | weak | • Full: $\tilde{O}(T^{1/2})$ <br> • Partial: $\tilde{O}(T^{3/4})$ | • Full: $\Omega(T^{1/2})$ <br> • Partial: $\Omega(T^{3/4})$ |
| **This paper** | adversarial setting | **global** | • Full: $\tilde{O}(T^{1/2})$ <br> • Partial: $\tilde{O}(T^{3/4})$ | • Full: $\Omega(T^{1/2})$ <br> • Partial: $\Omega(T^{5/7 \approx 0.714})$ |

Table 1: Comparison of prior results on bilateral trade. The positive result for a stochastic adversary in the partial feedback, marked with an asterisk ($*$), holds under the assumption that the seller and buyer valuations are drawn independently from smooth distributions. All the bounds in the second row (Azar et al. [2022]), marked with a dagger ($\dagger$), apply to 2-regret.

distribution, and the learner has either full feedback, or some stringent assumptions regarding the sequence of valuations are enforced. Specifically, in partial feedback regime, valuations have to be drawn i.i.d. from a smooth distribution, independently for the seller and the buyer [Cesa-Bianchi et al., 2021, Cesa-Bianchi et al., 2024]. If the learner is only required to enforce (step-wise) *weak budget balance* (*i.e.*, $p_t \le q_t$ for each $t$), then Azar et al. [2022] provide a learning algorithm achieving sublinear 2-regret when the sequence of valuation is generated by an oblivious adversary.[1] They also show that this result is tight: no algorithm can achieve sublinear $(2 - \varepsilon)$-regret in the adversarial case, for any constant $\varepsilon > 0$. In an attempt to overcome this barrier, Cesa-Bianchi et al. [2023] show that sublinear regret can be achieved beyond the i.i.d. stochastic setting, under the assumption that the adversary is constrained to choose randomized (possibly non-stationary) sequences of valuations that are not "too concentrated" (*i.e.,* under a $\sigma$-smooth adversary model). Inspired by the positive results obtained in the literature by transitioning from strong to weak budget balance, we investigate the following natural open question:

*Is it possible to achieve sublinear regret against an oblivious adversary in the repeated bilateral trade problem under a realistic notion of budget balance?*

We answer this question positively by introducing *global budget balance*, where the learner is required to maintain budget balance only "overall". The idea behind global budget balance is to allow the learner to reinvest the profit gained in previous rounds (obtained by posting a lower price for the seller compared to the buyer), with the constraint that the learner cannot subsidize the market *over the whole time horizon*. Formally, a learning algorithm that posts prices $(p_1, q_1), (p_2, q_2), \ldots$ is global budget balanced if the following inequality holds almost surely: $\sum_{t=1}^{T} \text{PROFIT}_t(p_t, q_t) \ge 0$. The profit $\text{PROFIT}_t(p_t, q_t) = \mathbb{I}\{s_t \le p_t\}\mathbb{I}\{q_t \le b_t\}(q_t - p_t)$ is non-negative when $p_t \le q_t$, and may drop below zero only by posting prices that are not step-wise budget balanced, *i.e.*, $p_t > q_t$. We argue that this constraint is more realistic than the restrictive notions of per-round budget balance. For instance, in contexts like ride-hailing platforms (such as Uber and Lyft), the platform might opt to forego some short-term profit to enhance other metrics, like the overall welfare of the system.

---

[1]The $\alpha$-regret measures the difference between the gain from trade of the best fixed price in hindsight and $\alpha$ times that of the algorithm (see e.g., Kakade et al. [2009]).

## 1.1 Overview of Our Results

We report here an overview of our results, we also refer to Table 1 for a comparison with the state of the art. In this paper we introduce the notion of global budget balance for the repeated bilateral trade problem, and provide the following results in terms of regret with respect to the best fixed price in hindsight in the adversarial case:

- In the full feedback model, when the learner observes seller and buyer valuations after posting prices, we design a learning algorithm characterized by a $\tilde{O}(T^{1/2})$ regret upper bound (Theorem 4.2). We also prove that no learning algorithm can improve this bound by more than a poly-log $T$ factor (Theorem 4.4).

- In the *one-bit feedback* model, where the learner can observe only whether the trade happened or not, we show that it is possible to guarantee a $\tilde{O}(T^{3/4})$ regret upper bound (Theorem 5.4). Then, we provide an $\Omega(T^{\frac{5}{7} \approx 0.714})$ lower bound, which holds even in the *two-bit feedback* model, where the learner can observe which agent accepted and who declined the offered prices (Theorem 5.5).

These results demonstrate how the notion of global budget balance enables online learnability, allowing us to provide the first no-regret algorithms for repeated bilateral trade within an oblivious adversary framework, in contrast to the per-round approaches considered in previous works. Furthermore, the regret rates separate full feedback and the two partial feedback models (one or two bits). In partial feedback, the surprising lower bound of $\Omega(T^{5/7})$, together with the $O(T^{3/4})$ upper bound, mark a clear separation between this problem and other partial feedback models (*e.g.,* partial monitoring [Bartók et al., 2014] and online learning with feedback graph [Alon et al., 2017], where the minimax regret have been characterized to fall in one of three admissible rates: $\sqrt{T}$, $T^{2/3}$ and $T$). This separation had already been hinted at in the special case of $\sigma$-smooth adversary by Cesa-Bianchi et al. [2023].

Finally, inspired by work on bandits with knapsacks (see Section 1.3 for detailed references), we introduce a stronger learning benchmark: the best fixed feasible distribution over prices. Such benchmark is allowed to post prices that are not per-round budget balanced, but is global budget balanced in "expectation".

- We show that there exists a constant $\varepsilon_0 > 0$ such that it is impossible to achieve sublinear $\alpha$-regret against this benchmark for any $\alpha \in [1, 1 + \varepsilon_0)$ (Theorem 6.2).

- We prove that the best feasible distribution over prices collects at most twice the gain from trade extracted by the best fixed price in hindsight (Theorem 6.3). This implies the existence of algorithms with sublinear 2-regret against this new benchmark.

- We show that the multiplicative gap of 2 between the gain from trade attainable by the two different benchmarks is tight (Theorem 6.5).

First, we observe that the task of learning the best feasible distribution over prices is reminiscent of the problem of bandits with knapsacks in the presence of replenishment [Kumar and Kleinberg, 2022, Slivkins et al., 2023, Bernasconi et al., 2024a]. In contrast to previous work, we consider the more challenging adversarial setting and provide learning algorithms with a competitive ratio that is an absolute constant. In the adversarial bandits with knapsacks literature, the only setting where sublinear $\Theta(1)$-regret can be achieved is when the available budget is $\Omega(T)$ [Castiglioni et al., 2022], while in general the competitive ratio is $O(\log T)$ [Immorlica et al., 2022]. Second, the tight multiplicative gap of 2 between the two benchmarks suggests that to design a better learning

algorithm with sublinear $\alpha$-regret with respect to the best feasible distribution (for $\alpha \in (1 + \varepsilon_0, 2)$), a more direct approach is needed.

## 1.2 Challenges and Techniques

The key aspects that distinguish bilateral trade from standard online learning models with full or bandit feedback can be identified in two main features: the action space and the challenging partial feedback structure. The applicability of previous results to our model is significantly limited due to adversarial input sequences and the need to handle the global budget balance constraint effectively.

**Action space.** The action space is continuous and bidimensional (prices belong to $[0, 1]^2$), and neither the gain from trade nor the profit functions are continuous in the prices posted. This makes it challenging to discretize the space with a finite grid $G$ such that the best prices in $G$ perform similarly to the best prices in $[0, 1]^2$, and such that grid $G$ is small enough that it is possible to learn in an online way its best pair of prices. In the absence of any probabilistic or smoothness assumption on the adversary, we cannot rely on a "smoothing trick" to induce regularity on the expected gain from trade, as in previous works [Cesa-Bianchi et al., 2023].

**Partial Feedback.** Partial feedback models for bilateral trade are inherently challenging. The one-bit feedback model only informs the learner on whether the trade happened or not, which is significantly less informative than the traditional bandit feedback model, since the learner cannot even reconstruct the gain from trade received for the specific prices it posted. For example, if the learner posts price $1/2$ to both agents, and they accept the trade, there is no way of distinguishing between the case in which the gain from trade is constant (*e.g.*, valuations are $(0, 1)$) from the case in which the gain from trade is arbitrarily small (*e.g.*, valuations are $(1/2 - \varepsilon, 1/2 + \varepsilon)$ for some small $\varepsilon$). On the other hand, if one of the two agents rejects the trade, then the learner can only infer loose bounds on the valuations.

**Gain from Trade vs. Profit trade-off.** Global budget balance requires that the cumulative sum of profits at the end of the time horizon must be greater than or equal to 0. Therefore, the learner has to maximize its cumulative gain from trade, while accumulating enough profit to enforce global budget balance. Balancing this trade-off is a complex task due to the different nature of the two objectives: gain from trade is maximized by setting identical prices for both agents, whereas profit is maximized by selecting prices that are "far from each other". To see this, consider an instance where valuations are either $(s_t, b_t) = (0, 1)$ or $(s_t, b_t) = (1/2 - \varepsilon, 1/2 + \varepsilon)$ with equal probability, for some small $\varepsilon > 0$. To achieve maximum expected profit, the learner would always set the price at 0 for the seller and 1 for the buyer. On the other hand, to maximize the expected gain from trade, the learner would always offer $1/2$ to both agents.

**Our Two-Phase Approach.** Our learning algorithms follow a two-phase approach, initially focusing on maximizing profit through a carefully designed multiplicative grid $F_K$ of candidate prices and then switching to maximizing gain from trade on a different (additive) grid $H_K$ of non-budget-balanced prices. At a high level, the first phase is used to collect budget, which can be subsequently reinvested in the second phase. This poses several challenges due to the non-stationary nature of the adversary. The pairs of prices in $H_K$, which are not per-round budget balanced, enable the algorithm to circumvent the negative results that hinder discretization in scenarios with per-round budget balance (see, *e.g.*, , the "needle in a haystack" phenomenon in Theorem 7 of Cesa-Bianchi et al. [2024]). The multiplicative nature of the grid $F_K$ is crucial in ensuring that the gain from trade accrued by the algorithm during the first phase does not yield too much regret. This last result is surprising since, in the first phase, the learning algorithm is maximizing profit, an objective that

4

is inherently orthogonal to the gain from trade. Finally, the scarcity of feedback in the one-bit feedback model is addressed via a carefully designed estimation technique that allows the learner to estimate the gain from trade in one point of the grid $H_K$ posting two different prices. In contrast to the technique by Azar et al. [2022], our procedure is "asymmetric" in how it deals with the seller and buyer, and it provides biased estimates.

**Lower bounds.** Besides the typical challenges in proving lower bounds for repeated bilateral trade with respect to the best fixed price in hindsight, in our model the agent is allowed to post prices that are not per-round budget balanced (*i.e.,* it may be the case that $p_t > q_t$). This considerably complicates the construction of the hard instances, as any algorithm could sacrifice temporarily some profit by posting prices with $p_t > q_t$ to extract a large gain from trade (that the fixed price benchmark may not be able to obtain). To deter this kind of behavior, we incorporate into the hard instances certain unfavorable trade opportunities that dissuade the learner from setting prices that are not budget balanced. This additional complication comes at some cost: in the partial (two-bit) feedback model we recover a lower bound of $\Omega(T^{5/7})$, whereas the corresponding lower bound by Cesa-Bianchi et al. [2023] is $\Omega(T^{3/4})$.

## 1.3   Further Related Works

**Online Learning and Economics.** Regret minimization techniques have found applications across different domains motivated by economics, with the goal of overcoming unrealistic assumptions. For example, they have been applied to one-sided pricing [Kleinberg and Leighton, 2003, Feldman et al., 2016], auctions [Morgenstern and Roughgarden, 2015, Cesa-Bianchi et al., 2015, Lykouris et al., 2016, Weed et al., 2016, Balseiro and Gur, 2019, Nedelec et al., 2022, Daskalakis and Syrgkanis, 2022, Cesa-Bianchi et al., 2024], contract design [Ho et al., 2016, Zhu et al., 2023, Duetting et al., 2023], brokerage [Bolić et al., 2024], and Bayesian persuasion [Castiglioni et al., 2020, Zu et al., 2021, Castiglioni et al., 2023, Bernasconi et al., 2023].

**Partial feedback.** Repeated bilateral trade naturally involves challenges due to partial feedback. Therefore, our work aligns with the research that explores online learning with feedback models beyond the conventional full feedback and bandit models. Our one- and two-bit feedback models share similarities with *graph-structured feedback* [Alon et al., 2017] and with the *partial monitoring* framework [Cesa-Bianchi et al., 2006, Bartók et al., 2014].

**Bandits with knapsacks.** Another related line of work is that of online learning under long-term constraints. Some works study the case of static constraints and develop projection-free algorithms with sublinear regret and constraint violations [Mahdavi et al., 2012, Jenatton et al., 2016], while others study the case of time-varying constraints [Mannor et al., 2009, Yu et al., 2017, Sun et al., 2017]. Badanidiyuru et al. [2018] introduced and solved the (stochastic) bandits with knapsacks (BwK) framework, in which they consider bandit feedback and stochastic objective and cost functions. In this model, the learner's objective is to maximize utility while guaranteeing that, for each of the $m$ available resources, cumulative costs are below a certain budget $B$. Other optimal algorithms for stochastic BwK were proposed by Agrawal and Devanur [2019], Immorlica et al. [2022]. The setting with adversarial inputs was first studied in Immorlica et al. [2022], where the baseline considered is the best fixed distribution over arms. Achieving no-regret is not possible under this baseline and, therefore, they provide no-$\alpha$-regret guarantees for their algorithm. If we denote by $\rho$ the per-iteration budget of the learner, the best-known guarantees on the competitive ratio $\alpha$ are $1/\rho$ in the case in which $B = \Omega(T)$ [Castiglioni et al., 2022], and $O(\log m \log T)$ in the general case [Kesselheim and Singla, 2020]. When considering a benchmark similar to the adversarial BwK

**Learning Protocol of Repeated Bilateral Trade**

---

1  Initial budget $B_0 = 0$
2  **for** $t = 1, 2, \dots$ **do**
3      The adversary privately chooses $(s_t, b_t)$ in $[0, 1]^2$
4      The learner posts prices $(p_t, q_t) \in [0, 1]^2$ such that $p_t - q_t \leq B_t$
5      The learner receives a (hidden) reward $\mathrm{GFT}_t(p_t, q_t) \in [-1, 1]$
6      The budget of the learner is updated $B_t \leftarrow B_{t-1} + \mathrm{PROFIT}_t(p_t, q_t)$
7      Feedback $z_t$ is revealed to the learner

---

scenario, we show that our algorithm ensures a $\alpha = 2$ guarantee. Kumar and Kleinberg [2022] recently proposed a generalization of the stochastic BwK model in which resource consumption can be non-monotonic; that is, resources can be replenished or renewed over time. Our model also admits replenishment. It should be noted that, in our setting, directly utilizing techniques from BwK is not feasible due to the complex continuous action space and the limited availability of feedback, which is less informative compared to traditional bandit feedback.

## 2 Repeated Bilateral Trade

We study repeated bilateral trade problem in an online learning setting, where the learner has to enforce global budget balance and the sequence of valuations is generated by an oblivious adversary.

**The learning protocol.** The learner repeatedly interacts with the environment according to the following protocol (see also pseudocode). At each time step $t$, a new pair of buyer and seller arrives, characterized by valuations $b_t \in [0, 1]$ and $s_t \in [0, 1]$, respectively. Without knowing $s_t$ and $b_t$, the learner posts two prices: $p_t \in [0, 1]$ to the seller, and $q_t \in [0, 1]$ to the buyer. If both the seller and the buyer accept (*i.e.*, $s_t \leq p_t$ and $q_t \leq b_t$), then the learner is awarded the gain from trade

$$\mathrm{GFT}_t(p_t, q_t) = \mathbb{I}\{s_t \leq p_t\}\mathbb{I}\{q_t \leq b_t\}(b_t - s_t),$$

that corresponds to the increase in social welfare generated by the trade. To simplify the notation, we omit the second argument of $\mathrm{GFT}_t$ (and of $\mathrm{PROFIT}_t$) when the same price is posted to both agents. After posting the prices, the learner does *not* observe directly the gain from trade or the valuations, but receives some feedback $z_t$.

**Global budget balance.** For each time step $t$, the notion of *profit* of the learner is naturally defined: if the agents accept prices $p_t$ and $q_t$, then the learner receives a net profit of $q_t - p_t \in [-1, 1]$. Unlike the case of the gain from trade, the learner naturally knows its profit at the end of each time step, as it sets the prices and always observes whether the trade occurred. The learner maintains a budget $B_t$, which is initially 0 ($B_0 = 0$) and is updated at each time step according to the profit generated or consumed: $B_t \leftarrow B_{t-1} + \mathrm{PROFIT}_t(p_t, q_t)$. We restrict the learner to enforce a *global budget balance* property which states that the final budget $B_T$ has to be non-negative with probability 1. In practice, we require the learner to always post prices $p_t, q_t$ such that $(p_t - q_t) \leq B_{t-1}$.[2]

**Feedback models.** In this paper, we study three feedback models, that we list here in increasing order of intricacy:

---

[2]In fact, this condition is not just sufficient, but also necessary. Indeed, if $p_t - q_t > B_t$, the adversary might select valuations $(s_t, b_t)$ such that $\mathrm{PROFIT}_t(p_t, q_t) < -B_{t-1}$ and thus $B_t < 0$. After that, the adversary might select valuations $(s_\tau, b_\tau) = (0, 0)$ for all $\tau \geq t + 1$, thereby forcing $B_T = B_t < 0$.

- *Full feedback*: at the end of each round, the agents reveal their valuations (*i.e.*, $z_t = (s_t, b_t)$).

- *Two-bit feedback*: the agents only reveal their willingness to accept the prices offered by the learner (*i.e.*, $z_t$ is composed by the two bits $(\mathbb{I}\{s_t \leq p_t\}, \mathbb{I}\{q_t \leq b_t\})$)

- *One-bit feedback*: the learner only observes whether the trade happened or not (*i.e.*, $z_t = \mathbb{I}\{s_t \leq p_t\} \cdot \mathbb{I}\{q_t \leq b_t\}$).

These feedback models are not only interesting from the theoretical learning perspective, but they are also well motivated in terms of practical applications. The full-feedback model can be used to describe sealed-bid-type auctions, while the two partial feedback settings (one- and two-bit) enforce the desirable property (for the agents) of revealing a minimal amount of information to the learner.

**Regret with respect to the best fixed price.** The goal is to maximize the total gain from trade on a fixed and known time horizon $T$ while enforcing the global budget balance condition. Following the literature on repeated bilateral trade [Cesa-Bianchi et al., 2021], we measure the performance of a learning algorithm in terms of its regret with respect to the best fixed price(s) in hindsight. For any learning algorithm $\mathcal{A}$ and sequence of valuations $\mathcal{S} = \{(s_t, b_t)\}_{t=1}^T$ we define:

$$R_T(\mathcal{A}, \mathcal{S}) = \max_{\substack{(p,q)\in[0,1]^2 \\ p \leq q}} \sum_{t=1}^T \text{GFT}_t(p, q) - \mathbb{E}\left[\sum_{t=1}^T \text{GFT}_t(p_t, q_t)\right], \tag{1}$$

where the sequence $\mathcal{S}$ induces the $\text{GFT}_t$ functions and the expectation is with respect to (possibly) randomized prices $p_t$ and $q_t$ generated by the learning algorithm $\mathcal{A}$. One simple property that follows immediately by definition is that, for any sequence of valuations, there exists a fixed pair of identical prices that maximizes the gain from trade. This means that the notion of "best price in hindsight" is well defined, and confirms the intuition that posting two different prices only helps during learning, but does not impact the maximization of gain from trade in hindsight. Finally, we define the regret of an algorithm $\mathcal{A}$ (without the dependence on a specific sequence of valuations) as its worst-case performance: $R_T(\mathcal{A}) = \sup_{\mathcal{S}} R_T(\mathcal{A}, \mathcal{S})$, where the sup is over the set of all the possible sequences of $T$ pairs of valuations.

**A stronger benchmark: the best feasible distribution over prices.** In this paper we also introduce a new (stronger) benchmark for the study of repeated bilateral trade: the best fixed budget-feasible distribution over prices. This benchmark captures the flexibility of the global budget balance condition, and it arises naturally from the literature on *bandits with knapsacks*. Before proceeding with the definition, let $\Delta([0,1]^2)$ be the family of all the probability measures over the measurable space $([0,1]^2, \mathcal{B}([0,1]^2))$, where $\mathcal{B}$ denotes the Borel $\sigma$-algebra.

**Definition 2.1** (Best feasible distribution). For any sequence $\mathcal{S}$ of seller's and buyer's valuations, we define the best fixed budget-feasible distribution over prices as the solution of:

$$\sup_{\gamma \in \Delta([0,1]^2)} \mathbb{E}_{(p,q)\sim\gamma}\left[\sum_{t=1}^T \text{GFT}_t(p, q)\right] \tag{2}$$

$$\text{s.t.} \quad \mathbb{E}_{(p,q)\sim\gamma}\left[\sum_{t=1}^T \text{PROFIT}_t(p, q)\right] \geq 0,$$

where $\mathbb{E}_{(p,q)\sim\gamma}$ denotes that the expectation is with respect to prices $(p, q)$ sampled according to $\gamma$.

This definition is well posed and there exist optimal distributions whose support contains either one or two pairs of prices. For a formal proof of this fact we refer to Proposition A.3 in Appendix A.1.

# 3  Price Discretizations and Two-Phase Algorithm

In this section we present our two-phase meta algorithm, preceeded by two key results on how to discretize the price space in a way that ensures certain essential properties about profit and gain from trade. First, in Section 3.1 we prove that the gain from trade of the best fixed price in hindsight is close to that of the best pair of (non-budget-balanced) prices on a suitable "additive" grid. Second, in Section 3.2 we construct an hybrid "multiplicative-additive" grid in which each interval of a one-dimensional additive grid is further divided into sub-intervals with geometrically decreasing length. This grid has the surprising property that the profit of the best fixed pair of prices on it is close to the gain from trade generated by the best fixed price in the $[0, 1]$ interval, up to a poly-logarithmic multiplicative factor. Finally, we introduce our two-phase learning via the meta-algorithm GFT-Max.

## 3.1  Additive Grid for Gain from Trade

For any integer $K$, we denote by $G_K = \{0, 1/K, 2/K, \ldots, 1 - 1/K, 1\}$ the *K-uniform grid* over $[0, 1]$. Similarly, we denote with $H_K = \{(i+1/K, i/K) : i \in \{0, 1, \ldots, K - 1\}\}$ the set of pairs formed by contiguous points in the $K$-uniform grid such that the first element of the pair is greater than the second. This latter grid can be proved to enjoy the desirable property of well-approximating the gain from trade of the best fixed price, while violating the global budget balance condition by a small amount. The argument behind the approximation guarantee is simple: if $p^*$ is the best fixed price in hindsight, then the pair of prices $((i+1)/K, i/K)$ such that $p^*$ belongs to the interval $[i/K, (i+1)/K]$ are nearly as good as $p^*$. We have the following result.

**Proposition 3.1.** *For any $K$ and sequence of valuations, we have:*

$$\max_{p \in [0,1]} \sum_{t=1}^{T} GFT_t(p) \le \max_{(p,q) \in H_K} \sum_{t=1}^{T} GFT_t(p, q) + \frac{T}{K}.$$

*For any $(p, q) \in H_K$, total profit $\sum_{t=1}^{T} \text{Profit}_t(p, q)$ is at least $-T/K$.*

*Proof.* The optimal price in hindsight $p^*$ is contained in some interval $[i^*/K, (i^*+1)/K]$. For any time $t$ we have the following cases:

   (*i*) If $GFT_t(p^*) = 0$, then the gain from trade of $((i^*+1)/K, i^*/K)$ is at least $-1/K$ (when the valuation of the seller is $(i^*+1)/K$, and that of the buyer is $i^*/K$).

   (*ii*) If $GFT_t(p^*) > 0$, then posting the pair of prices $((i^*+1)/K, i^*/K)$ makes the trade happen, and guarantees the same gain from trade.

Then, by summing up the gain from trade obtained by posting $((i^*+1)/K, i^*/K)$, we immediately obtain the first part of the statement by applying at each $t$ either case (*i*) or (*ii*). The second part of the statement follows from the observation that the per-round deficit for posting prices $((i^*+1)/K, i^*/K)$ is at most $1/K$. This concludes the proof. $\qquad\square$

A simple calculation shows that $GFT_t((i+1)/K, i/K)$ is bounded by the sum of $GFT_t(i/K)$ and $GFT_t((i+1)/K)$. Therefore, we obtain the following known result as a Corollary to Proposition 3.1.

**Corollary 3.2** (Claim 1 of Azar et al. [2022])**.** *For any $K$ and sequence of valuations, we have:*

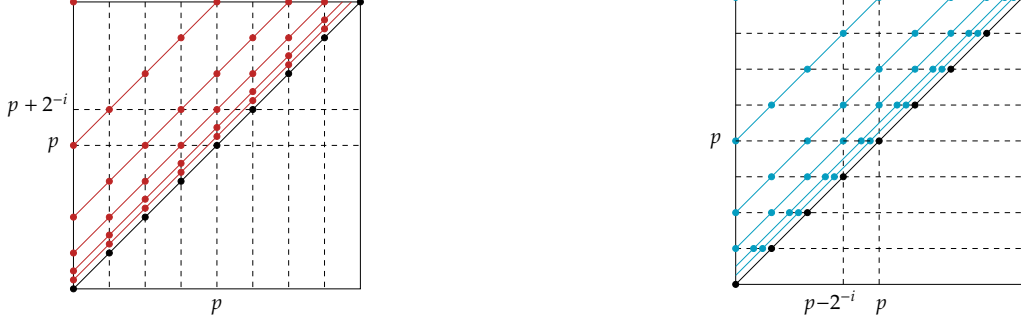$$\max_{p \in [0,1]} \sum_{t=1}^{T} GFT_t(p) \le 2 \cdot \max_{p \in G_K} \sum_{t=1}^{T} GFT_t(p) + \frac{T}{K}.$$

Figure 1: $F_K^+$ (left) and $F_K^-$ (right) for $K = 8$, $T = 32$, $i = 3$.

## 3.2 Multiplicative Grid for Profit

For any $K$, we construct the two-dimensional grid $F_K$ starting from the points on the one-dimensional grid $G_K$. For each $p \in G_K$, we add to $F_K$ points of the form $(p - 2^{-i}, p)$ and $(p, p + 2^{-i})$, for $i = 0, 1, \ldots, \lceil \log T \rceil$ so that they define intervals of geometrically decreasing length to the left and upward of $(p, p)$. Formally, we define $F_K$ as the union of $F_K^-$ and $F_K^+$ (see also Figure 1):

$$F_K^- = \left\{ (p - 2^{-i}, p) : p \in G_K \text{ and } i \in \{0, 1, \ldots, \lceil \log T \rceil\} \right\} \cap [0, 1]^2,$$
$$F_K^+ = \left\{ (p, p + 2^{-i}) : p \in G_K \text{ and } i \in \{0, 1, \ldots, \lceil \log T \rceil\} \right\} \cap [0, 1]^2.$$

The additive-multiplicative nature of $F_K$ endows it with two crucial properties: (i) its cardinality is $O(K \log T)$ an thus only depends linearly in $K$ and (ii) the profit of the best prices in $F_K$ is at least a $O(\log T)$ fraction of the GFT at the best fixed price in $[0, 1]$, up to an additive factor of $O(T/K)$.

**Proposition 3.3.** *For any K and sequence of valuations, we have:*

$$\max_{p \in [0,1]} \sum_{t=1}^{T} GFT_t(p) \leq 12 \log T \cdot \max_{(p,q) \in F_K} \sum_{t=1}^{T} PROFIT_t(p, q) + \frac{5T}{K}.$$

*Proof.* Fix the sequence $\mathcal{S}$ of valuations and let $p^*$ be the price maximizing the gain from trade in $G_K$. We have the following chain of inequalities:

$$\max_{p \in [0,1]} \sum_{t=1}^{T} GFT_t(p) \leq 2 \sum_{t=1}^{T} (b_t - s_t) \mathbb{I}\{s_t \leq p^* \leq b_t\} + \frac{T}{K} \qquad \text{(by Corollary 3.2)}$$

$$= 2 \sum_{t=1}^{T} (b_t - p^*) \mathbb{I}\{s_t \leq p^* \leq b_t\} + 2 \sum_{t=1}^{T} (p^* - s_t) \mathbb{I}\{s_t \leq p^* \leq b_t\} + \frac{T}{K}. \qquad (3)$$

We bound separately the first and second term of the right-hand side of the inequality. Starting with $(b_t - p^*) \mathbb{I}\{s_t \leq p^* \leq b_t\}$, we can rewrite the expression through a case analysis depending on the interval of the discretization in which $b_t$ is located. For each time step $t$, we have

$$(b_t - p^*) \mathbb{I}\{s_t \leq p^* \leq b_t\} \leq \sum_{i=0}^{\lceil \log T \rceil} (b_t - p^*) \mathbb{I}\{s_t \leq p^*\} \mathbb{I}\{p^* + 2^{-i} \leq b_t < p^* + 2^{-i+1}\} + \frac{1}{T},$$

9

where we used the fact that $b_t - p^* \leq 1/T$ if $b_t$ belongs to $[p_t^*, p_t^* + 1/T]$. This yields

$$(b_t - p^*)\mathbb{I}\{s_t \leq p^* \leq b_t\} \leq 2 \cdot \sum_{i=0}^{\lceil \log T \rceil} 2^{-i}\mathbb{I}\{s_t \leq p^*\}\mathbb{I}\{p^* + 2^{-i} \leq b_t < p^* + 2^{-i+1}\} + \frac{1}{T}. \quad (4)$$

Let $n_i$ be the number of time steps satisfying the condition $\{s_t \leq p^*, p^* + 2^{-i} \leq b_t < p^* + 2^{-i+1}\}$. Summing up Equation (4) for $t = 1, 2, \ldots T$ we get

$$\sum_{t=1}^{T}(b_t - p^*)\mathbb{I}\{s_t \leq p^* \leq b_t\} \leq 2 \cdot \sum_{i=0}^{\lceil \log T \rceil} \frac{n_i}{2^i} + \frac{T}{T} \leq 3\log T \cdot \max_{i \in \{0,\ldots,\log T\}} \frac{n_i}{2^i} + 1$$

$$\leq 3\log T \cdot \max_{(p,q) \in F_K} \sum_{t=1}^{T} \mathrm{P}\mathrm{R}\mathrm{O}\mathrm{F}\mathrm{I}\mathrm{T}_t(p, q) + 1. \quad (5)$$

To obtain Equation (5) we use that, for any $i \in \{0, \ldots, \lceil \log T \rceil\}$, if $n_i > 0$ then it must be the case that $p^* + 2^{-i} \in [0, 1]$. Therefore, for any $i$, it it possible to obtain a profit of $2^{-i}$ by posting the pair $(p^*, p^* + 2^{-i})$, which is guaranteed to belong to $F_K$ since $p^* \in G_K$ by construction, and $p^* + 2^{-i} \in [0, 1]$. A similar argument can be carried over for the other term of Equation (3), yielding:

$$\sum_{t=1}^{T}(p^* - s_t)\mathbb{I}\{s_t \leq p^* \leq b_t\} \leq 3\log T \cdot \max_{(p,q) \in F_K} \sum_{t=1}^{T} \mathrm{P}\mathrm{R}\mathrm{O}\mathrm{F}\mathrm{I}\mathrm{T}_t(p, q) + 1. \quad (6)$$

Finally, we plug Equations (5) and (6) into Equation (3), and use $K \leq T$ to conclude the proof. □

### 3.3 Our Two-Phase Meta-Algorithm: GFT-Max

We describe our two-phase learning approach by presenting the meta-algorithm GFT-Max. For details we refer to the pseudocode. The algorithm takes in input a budget threshold $\beta$ and an integer $K$ (which induces the two grids $F_K$ and $H_K$), and employs two regret minimizers—$\mathcal{A}_P$ for the profit and $\mathcal{A}_G$ for the gain from trade—as internal routines. In the first phase (Line 1), the algorithm uses function Profit-Max to maximize profit until the collected budget reaches a given threshold $\beta$. This is achieved by running a regret minimizer $\mathcal{A}_P$ over the set $F_K$ of pairs of prices (see Section 3.2) using profit as objective. Then, in the second phase (from Line 2 onward), the algorithm exploits a regret minimizer $\mathcal{A}_G$ to maximize the gain from trade over the grid $H_K$, whose prices which are "almost budget-balanced" and consume only a small fraction of the previously acquired budget (see Proposition 3.1). In Section 4 and Section 5 we provide regret upper bounds for this meta-algorithm in the full and one-bit feedback model, respectively. The budget threshold $\beta$, the regret minimizers, and the grid parameter $K$ are tuned according to the specific case considered.

## 4 Full Feedback

We start by studying the *full feedback* input model where the agents reveal their valuations $(s_t, b_t)$ at the end of each time step $t$. Here, the learner has counterfactual information regarding all the prices they could have posted, *independently* of the pair of prices actually posted at time $t$. In Section 4.1, we first present a two-phase learning algorithm (GFT-Max) which guarantees $\tilde{O}(\sqrt{T})$ regret with respect to the best fixed price in hindsight. In Section 4.2 we complement this result by proving that this is tight, up to poly-logarithmic terms.

**Algorithm 1:** GFT-Max

> **Input:**  • budget threshold $\beta$
> • integer $K$ and price-grids $F_K$ and $H_K$
> • regret minimizers $\mathcal{A}_P$ and $\mathcal{A}_G$

1 Run Profit-Max $(\beta, F_K, \mathcal{A}_P)$             /* Phase I */
2 **if** *Profit-Max terminated at time step $\tau < T$* **then**
3     Initialize $\mathcal{A}_G$ on $H_K$             /* Phase II */
4     **for** $t = \tau + 1, 2, \ldots, T$ **do**
5        Receive from $\mathcal{A}_G$ the prices $(p_t, q_t)$
6        Post prices $(p_t, q_t)$ and observe feedback $z_t$
7        Feed feedback $z_t$ to $\mathcal{A}_G$

8 **function** *Profit-Max $(\beta, F_K, \mathcal{A}_P)$*
> **Input:**  • budget threshold $\beta$
> • grid $F_K$ of pairs of prices
> • regret minimizer $\mathcal{A}_P$

9     Initialize $\mathcal{A}_P$ on $|F_K|$ actions, one for each $(\hat{p}, \hat{q}) \in F_K$, and set $B_0 \leftarrow 0$
10     **for** $t = 1, 2, \ldots, T$ **do**
11        Receive from $\mathcal{A}_P$ the prices $(p_t, q_t)$
12        Post prices $(p_t, q_t)$ and observe feedback $z_t$
13        Feed feedback $z_t$ to $\mathcal{A}_P$
14        Update $B_t \leftarrow B_{t-1} + \text{Profit}_t(p_t, q_t)$
15        **if** $B_t \geq \beta$ **then** Terminate the algorithm

## 4.1 $\tilde{O}(\sqrt{T})$ Upper Bound with Full Feedback

We start the analysis by looking at the first phase of GFT-Max, Profit-Max (reported as a function in the pseudocode of GFT-Max). We employ the Hedge algorithm (see, *e.g.*, Section 5.3 of Slivkins [2019]) as the regret minimizer $\mathcal{A}_P$, which is used on the action space of the prices in $F_K$. As a first step, we note that the gain from trade of any fixed price in the first phase (which terminates at the stopping time $\tau$) is not too large.

**Lemma 4.1.** *Consider Profit-Max with budget threshold $\beta$, grid $F_K$, and learning algorithm Hedge as $\mathcal{A}_P$. Then, with probability at least $1 - 1/T$, we have*

$$\max_{p \in [0,1]} \sum_{t=1}^{\tau} GFT_t(p) \leq 8(\beta + 1) \log T + \frac{5T}{K} + 32 \log T \sqrt{T \log(T|F_K|)}.$$

*Proof.* We start by observing that, by Proposition 3.3, there exists a pair of prices $(p^*, q^*) \in F_K$ such that

$$\max_{p \in [0,1]} \sum_{t=1}^{\tau} \text{GFT}_t(p) \leq 8 \log T \cdot \sum_{t=1}^{\tau} \text{Profit}_t(p^*, q^*) + \frac{5T}{K}.$$

Hedge maintains a distribution $\gamma_t \in \Delta(F_K)$ at each $t \in [T]$, and such distributions guarantees that the expected regret is $O(\sqrt{T \log(|F_K|)})$ [Slivkins, 2019]. In particular, given $s \in [T]$, we have

$$\sum_{t=1}^{s} \text{Profit}_t(p^*, q^*) \leq \sum_{t=1}^{s} \mathop{\mathbb{E}}_{(p,q) \sim \gamma_t} [\text{Profit}_t(p, q)] + 2\sqrt{T \log(|F_K|)}.$$

11

By applying the Azuma-Hoeffding inequality for each round $s \in [T]$, and union bounding over the possible stopping times, we get that with probability at least $1 - 1/T$, we can write the following also for the stopping time $\tau$:

$$\sum_{t=1}^{\tau} \text{PROFIT}_t(p^*, q^*) \leq \sum_{t=1}^{\tau} \text{PROFIT}_t(p_t, q_t) + 2\sqrt{T \log(|F_K|)} + 4\sqrt{T \log(T)}.$$

This yields the following chain of inequalities

$$\max_{p \in [0,1]} \sum_{t=1}^{\tau} \text{GFT}_t(p) \leq 8 \log T \cdot \sum_{t=1}^{\tau} \text{PROFIT}_t(p^*, q^*) + \frac{5T}{K}$$

$$\leq 8 \log T \cdot \sum_{t=1}^{\tau} \text{PROFIT}_t(p_t, q_t) + \frac{5T}{K} + 16 \log T \sqrt{T \log(|F_K|)} + 32 \log T \sqrt{T \log(T)}$$

$$\leq 8 \log T \cdot B_\tau + \frac{5T}{K} + 32 \log T \sqrt{T \log(T|F_K|)}$$

$$\leq 8 \log T \cdot (\beta + 1) + \frac{5T}{K} + + \frac{5T}{K} + 32 \log T \sqrt{T \log(T|F_K|)}$$

This concludes the proof. □

Lemma 4.1 helps us bounding the regret of GFT-MAX up to the (random) time step $\tau$, when the algorithm switches from profit to gain from trade maximization. Setting $\beta = \sqrt{T}$ and $K = \sqrt{T}$, and using HEDGE as regret minimizer also in the second phase, yields the following result.

**Theorem 4.2.** *Consider the repeated bilateral trade problem in the full feedback model. There exists a learning algorithm $\mathcal{A}$ that respects global budget balance and whose regret with respect to the best fixed price in hindsight verifies*

$$R_T(\mathcal{A}) \leq 92 \log^{3/2}(T) \sqrt{T}.$$

*Proof.* We prove that algorithm GFT-MAX with the proper choice of budget $\beta$, grids $F_K$ and $H_K$, and algorithms $\mathcal{A}_P$ and $\mathcal{A}_G$ achieves the desired regret bound, while enforcing global budget balance. First, we show that the algorithm enforces global budget balance for any value of the stopping time $\tau \in [T]$. By construction, the profit at time $\tau$ (*i.e.,* right after the end of first phase in which we employ the subroutine PROFIT-MAX) is at least $\beta$. Moreover, in each round $t \in \{\tau + 1, \ldots, T\}$ of the second phase, the profit is at least $-1/K$. Hence, the cumulative profit at time $T$ is at least

$$\beta - (T - \tau)\frac{1}{K} \geq \sqrt{T} - T\frac{1}{\sqrt{T}} = 0.$$

Then, we prove the upper bound on the cumulative regret. We start by considering the regret accumulated in the interval $\{\tau + 1, \ldots, T\}$. In particular, for any $\tau \in [T]$, we have

$$\max_{p \in [0,1]} \sum_{t=\tau+1}^{T} \text{GFT}_t(p) \leq \max_{(p,q) \in H_K} \sum_{t=\tau+1}^{T} \text{GFT}_t(p, q) + \frac{T}{K}$$

$$\leq \mathop{\mathbb{E}}_{(p_t, q_t) \sim \gamma_t} \left[ \sum_{t=\tau+1}^{T} \text{GFT}_t(p_t, q_t) \right] + \frac{T}{K} + 4\sqrt{T \log(|H_K|)}, \qquad (7)$$

12

where the first inequality follows from Proposition 3.1, and the second inequality follows from the regret bound of HEDGE when the range of the rewards is $[-1/K, 1]$ and $\gamma_t$ is the probability distribution over the action set maintained by HEDGE (instantiated to maximize gain from trade in the second phase).

Then, assume that the bound in Lemma 4.1 holds, which happens with probability at least $1 - 1/T$. By employing Equation (7) and Lemma 4.1 we can show that, with probability at least $1 - 1/T$,

$$
\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \leq \max_{p \in [0,1]} \sum_{t=1}^{\tau} \mathrm{GFT}_t(p) + \max_{p \in [0,1]} \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p)
$$

$$
\leq 8 \log T(\beta + 1) + \frac{5T}{K} + 32 \log T \sqrt{T \log(T|F_K|)} + \max_{p \in [0,1]} \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p)
$$

$$
\leq \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] + 8 \log T(\beta + 1) + \frac{6T}{K} + 32 \log T \sqrt{T \log(T|F_K|)} + 4\sqrt{T \log(|H_K|)},
$$

where the second inequality follows from Lemma 4.1, and the third one from Equation (7). Then, by substituting $\beta = K = \sqrt{T}$ we can conclude that

$$
\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \leq \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] + 8 \log T(\sqrt{T} + 1) + 6\sqrt{T}
$$

$$
+ 32 \log T \sqrt{T \log\left(2T^{3/2}(\log T + 1)\right)} + 4\sqrt{T \log(\sqrt{T})}
$$

$$
\leq \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] + 90 \log^{3/2}(T)\sqrt{T},
$$

By rearranging we have that, with probability at least $1 - 1/T$, it holds

$$
\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) - \mathbb{E}\left[ \sum_{t=1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] \leq \max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) - \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right]
$$

$$
\leq 90 \log^{3/2}(T)\sqrt{T},
$$

where the first inequality follows from the fact that the gain from trade is always non-negative. Finally, we can conclude that the expected regret is at most

$$
R_T(\text{GFT-MAX}) \leq \left(1 - \frac{1}{T}\right)\left(90 \log^{\frac{3}{2}}(T)\sqrt{T}\right) + \frac{1}{T} \cdot 2T \leq 92 \log^{\frac{3}{2}}(T)\sqrt{T}.
$$

This concludes the proof. □

## 4.2 $\Omega(\sqrt{T})$ Lower Bound with Full Feedback

We present a lower bound that shows how the regret rate in Theorem 4.2 is optimal up to poly-logarithmic factors. The lower bound is based on the following stochastic sequence: at each time step $t$ the pair $(s_t, b_t)$ is drawn uniformly at random between 3 pairs of valuations: $(0, 1/4)$, $(3/4, 1)$ and $(3/4, 1/4)$. These three points naturally partition the $[0, 1]^2$ square into four regions (see Figure 2). Crucially, prices in the $[3/4, 1] \times [0, 1/3]$ region (green in Figure 2) incur in negative expected gain
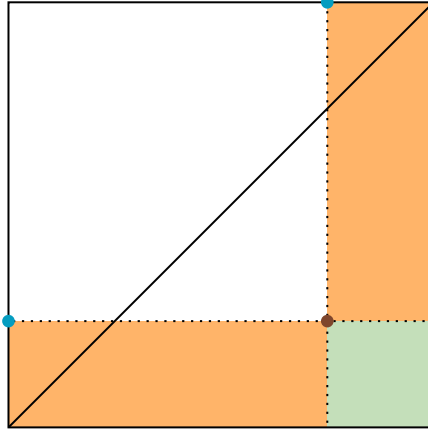
13

Figure 2: Partition of $[0, 1]^2$ as in the proof of Theorem 4.4.

from trade, while prices in the $[0, 3/4) \times (1/3, 1]$ region (white in Figure 2) miss all trades. Therefore, the only reasonable option for any learner is to post prices in the two remaining regions (**orange** in Figure 2), with an expected gain from trade of $1/12$. This allows for a reduction to an expert problem with 2 available actions (one for each of the two orange regions). This construction highlights a key difficulty if compared to lower bounds for per-round budget balanced algorithms: we need to disincentivize the learner from choosing non budget balanced prices below the diagonal. We have the following Theorem, which is preceded by a preliminary Lemma.

**Lemma 4.3.** *Let $S_n$ be a symmetric random walk on the line after $n$ steps, starting from $0$. Then, for n large enough, it holds that $\mathbb{E}[|S_n|] \geq 2/3\sqrt{n}$.*

*Proof.* It is well known that the expected distance of a random walk from the origin grows like $\Theta(\sqrt{n})$. Formally, the following asymptotic result holds (see, *e.g.,* Palacios [2008]):

$$\lim_{n \to \infty} \frac{\mathbb{E}[|S_n|]}{\sqrt{n}} = \sqrt{\frac{2}{\pi}}.$$

Observe that $\sqrt{2/\pi} > 2/3$. Then, there exists a finite $n_0$ such that $\mathbb{E}[|S_n|] \geq 2/3\sqrt{n}$ for all $n \geq n_0$. □

**Theorem 4.4.** *Consider the repeated bilateral trade problem in the full feedback model. Any learning algorithm that satisfies global budget balance suffers at least $\Omega(\sqrt{T})$ regret with respect to the best fixed price in hindsight.*

*Proof.* We prove this result via Yao's principle [Yao, 1977]. We apply the easy direction of the theorem, which reads (using our terminology) as follows: the regret $R_T(\mathcal{A})$ of a randomized learner $\mathcal{A}$ against the worst-case valuations sequence is at least the regret of the optimal deterministic learner $A$ against a stochastic sequence of valuations $\mathcal{S}$. Formally,

$$R_T(\mathcal{A}) \geq \sup_A \mathbb{E}\left[\max_{p \in [0,1]} \sum_{t=1}^{T} \text{GFT}_t(p) - \sum_{t=1}^{T} \text{GFT}_t(p_t, q_t)\right],$$

where the expectation is with respect to the stochastic valuation sequence $\mathcal{S}$, while $A$ denotes deterministic learner that posts the $(p_t, q_t)$ prices. In particular, we construct a randomized instance

$\mathcal{S}$ such that any deterministic learning algorithm must suffer, in expectation with respect to the randomness of $\mathcal{S}$, at least $c\sqrt{T}$ regret for some constant $c$.

The randomized instance is constructed as follows: at each time step $t \in [T]$ the adversary selects uniformly and independently at random one of the following three points $(0, 1/4)$, $(3/4, 1)$ and $(3/4, 1/4)$. We first compute a lower bound on the expected gain from trade achieved by the best fixed price in hindsight, and then we provide an upper bound on the expected gain from trade which can be attained by any deterministic learning algorithm. Combining these two intermediate results will yield the statement via Yao's principle.

Let $N_0$ be a random variable denoting the number of times that $(3/4, 1/4)$ is realized. Analogously, let $N_1$ (resp., $N_2$), be the number of times in which $(0, 1/4)$ (resp., $(3/4, 1)$) is realized. Clearly, $N_0 + N_1 + N_2 = T$, and $\mathbb{E}[N_i] = T/3$ for any $i = 0, 1, 2$. Conditioning on $N_0$, the remaining $T - N_0$ valuations are either $(0, 1/4)$ or $(3/4, 1)$, sampled uniformly and independently at random. Then, we have that

$$\mathbb{E}\left[\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \,\middle|\, N_0\right] \geq \mathbb{E}\left[\max_{p \in \{1/4, 3/4\}} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \,\middle|\, N_0\right]$$

$$= \mathbb{E}\left[\max_{p \in \{1/4, 3/4\}} \sum_{t=1}^{T} \frac{1}{4}\mathbb{I}\{s_t \leq p \leq b_t\} \,\middle|\, N_0\right]$$

$$= \frac{1}{4}\mathbb{E}[\max\{N_1, T - N_0 - N_1\} \mid N_0] \qquad \text{(Definitions of } N_i\text{)}$$

$$= \frac{1}{8}(T - N_0) + \frac{1}{8}\mathbb{E}[\max\{2N_1 - T + N_0, T - N_0 - 2N_1\} \mid N_0]$$

$$= \frac{1}{8}(T - N_0) + \frac{1}{8}\mathbb{E}[\max\{N_1 - N_2, N_2 - N_1\} \mid N_0]$$

$$= \frac{1}{8}(T - N_0) + \frac{1}{8}\mathbb{E}[|S_{T-N_0}| \mid N_0] \qquad (8)$$

$$\geq \frac{1}{8}(T - N_0) + \frac{\sqrt{T - N_0}}{12}, \qquad \text{(Lemma 4.3)}$$

where Equation (8) follows by considering a symmetric random walk on a line on $T - N_0$ steps that goes left when $(s_t, b_t) = (0, 1/4)$, and goes right when $(s_t, b_t) = (3/4, 1)$. Now, we can take the expectation (with respect to $N_0$) on the first and last term of the previous chain of inequalities to get

$$\mathbb{E}\left[\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p)\right] = \mathbb{E}\left[\mathbb{E}\left[\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \,\middle|\, N_0\right]\right] \qquad \text{(Conditional expectation)}$$

$$\geq \frac{1}{12}T + \frac{1}{12}\mathbb{E}\left[\sqrt{T - N_0}\right] \qquad (\mathbb{E}[N_0] = T/3)$$

$$\geq \frac{1}{12}T + \frac{\sqrt{T}}{24}\mathbb{P}(N_0 \leq 3/4T) \qquad \text{(Conditioning on } N_0 \leq 3/4T\text{)}$$

$$\geq \frac{1}{12}T + \frac{5\sqrt{T}}{216}, \qquad (9)$$

where the last line follows from Markov's inequality.

Now, we construct an upper bound on the gain from trade achievable by *any* deterministic learning algorithm (even without the constraint of enforcing global budget balance). Consider what happens at each fixed time steps $t$: the history of the realized valuations up to that point induce

15

deterministically the pair of prices $(p_t, q_t)$ posted by the learning algorithm. We prove now that no matter $(p_t, q_t)$ chosen, the learner does not achieve more than an expected gain from trade of $1/12$. To see this we study separately four cases:

- If $(p_t, q_t) \in [0, 3/4) \times (1/4, 1]$, then $\mathrm{GFT}_t(p_t, q_t) = 0$ with probability 1 because it misses all the possible trades.

- If $(p_t, q_t) \in [0, 3/4) \times [0, 1/4]$, then the learner gets $1/4$ gain from trade only when $(s_t, b_t) = (0, 1/4)$ is realized and 0 otherwise, for an expected gain from trade of $1/12$

- Similarly, if $(p_t, q_t) \in (3/4, 1] \times (1/4, 1]$, then the learner gets $1/4$ gain from trade only when $(s_t, b_t) = (3/4, 1)$ is realized and 0 otherwise, for an expected gain from trade of $1/12$

- Finally. if $(p_t, q_t) \in [3/4, 1] \times [0, 1/4]$, then the learner always observes a trade, but the expected gain from trade it achieves is 0 ($1/4$ with probability $2/3$ and $-1/2$ with the remaining probability).

Therefore, no matter what the learner does, it gets expected gain from trade at most $T/12$:

$$\mathbb{E}\left[\sum_{t=1}^{T} \mathrm{GFT}_t(p_t, q_t)\right] \leq \frac{T}{12}. \tag{10}$$

We can conclude the proof of the Theorem by combining Equation (9) and Equation (10) to get:

$$\mathbb{E}\left[\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) - \sum_{t=1}^{T} \mathrm{GFT}_t(p_t, q_t)\right] \geq \left(\frac{1}{12}T + \frac{5\sqrt{T}}{216}\right) - \frac{T}{12} = \frac{5\sqrt{T}}{216},$$

where that the randomness is with respect to the sequence generated by the randomized adversary. This concludes the proof. □

# 5 Partial Feedback

In this section, we study the more challenging partial feedback models. In Section 5.1, we provide a positive result for the case of one-bit feedback ($z_t = \mathbb{I}\{s_t \leq p_t\} \cdot \mathbb{I}\{q_t \leq b_t\}$), where the learner only observes whether the trade happened or not. In particular, we show that GFT-MAX, with a suitable initialization, achieves a regret of the order $\tilde{O}(T^{3/4})$. Differently from the full-information setting, the design of a no-regret algorithm for the gain from trade (*i.e.*, $\mathcal{A}_G$) is particularly challenging as we need to build an estimator for the gain from trade by only playing non-budget balanced prices in $H_K$.

In Section 5.2 we complement the regret upper bound by proving that every algorithm has regret at least $\Omega(T^{5/7})$, even with two-bit feedback ($z_t = (\mathbb{I}\{s_t \leq p_t\}, \mathbb{I}\{q_t \leq b_t\})$), *i.e.*, where each agent separately reveal their willingness to accept the prices posted. One of the main challenges posed by such a lower bound resides in handling non-budget balanced prices, as any algorithm could temporarily sacrifice some profit while collecting large GFT.

## 5.1 $\tilde{O}(T^{3/4})$ Upper Bound with One-Bit Feedback

We show how to employ GFT-MAX with a suitable choice of parameters $\beta$ and $K$, and regret minimizers $\mathcal{A}_P$ and $\mathcal{A}_G$ to achieve the desired regret bound. Section 5.1.1 presents a regret-minimizing algorithm that can be employed as $\mathcal{A}_P$, while Section 5.1.2 provides a suitable regret minimizer to be employed as $\mathcal{A}_G$. Finally, in Section 5.1.3, we present the final regret upper bound.

### 5.1.1 Regret Minimizer for Profit under Partial Feedback

As in the full-information setting, we exploit PROFIT-MAX to maximize the profit until the accrued budget is at least a given threshold $\beta$. In particular, we instantiate the subroutine PROFIT-MAX with EXP3.P [Auer et al., 2002] as regret minimizer $\mathcal{A}_P$ and grid $F_K$. The following lemma shows that the gain from trade of any fixed price $p$ in the first phase is small enough up to the stopping time $\tau$ that terminates the first phase.

**Lemma 5.1.** *Consider PROFIT-MAX with budget threshold $\beta$, grid $F_K$, and learning algorithm EXP3.P as $\mathcal{A}_P$. Then with probability at least $1 - 1/T$, we have that $\max_{p \in [0,1]} \sum_{t=1}^{\tau} GFT_t(p)$ is at most $8(\beta+1) \log T + \frac{5T}{K} + 256 \log T \sqrt{|F_K| T \log(|F_K|T)}$.*

*Proof.* First, note that by Proposition 3.3 there exists a pair of prices $(p^*, q^*) \in F_K$ such that

$$\max_{p \in [0,1]} \sum_{t=1}^{\tau} \text{GFT}_t(p) \leq 8 \log T \cdot \sum_{t=1}^{\tau} \text{PROFIT}_t(p^*, q^*) + \frac{5T}{K}.$$

Moreover, EXP.P guarantees that, with probability at least $1 - 1/T$, it holds

$$\sum_{t=1}^{\tau} \text{PROFIT}_t(p^*, q^*) \leq \sum_{t=1}^{\tau} \text{PROFIT}_t(p_t, q_t) + 32 \sqrt{|F_K| T \log(|F_K|T)}.$$

Then,

$$\max_{p \in [0,1]} \sum_{t=1}^{\tau} \text{GFT}_t(p) \leq 8 \log T \cdot \sum_{t=1}^{\tau} \text{PROFIT}_t(p^*, q^*) + \frac{5T}{K}$$

$$\leq 8 \log T \cdot \sum_{t=1}^{\tau} \text{PROFIT}_t(p_t, q_t) + \frac{5T}{K} + 256 \log T \sqrt{|F_K| T \log(|F_K|T)}$$

$$\leq 8 \log T B_\tau + \frac{5T}{K} + 256 \log T \sqrt{|F_K| T \log(|F_K|T)}$$

$$\leq 8 \log T (\beta + 1) + \frac{5T}{K} + 256 \log T \sqrt{|F_K| T \log(|F_K|T)}.$$

This concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

### 5.1.2 Regret Minimizer for Gain from Trade under Partial Feedback

A crucial ingredient we need is an estimation procedure capable of extracting quantitative information from the gain from trade, having only access to one bit of feedback. More precisely, we need an estimation procedure of the gain from trade function $H_K \ni (p, q) \mapsto \text{GFT}_t(p, q)$. A similar challenge is faced in Azar et al. [2022], where the action set consists of a discretization of a single price (*i.e.*, their estimation procedure posts $p$ to both seller and buyer). However, in our scenario, such symmetry no longer applies. Here, we must consider the grid $H_K$, which employs distinct prices for the seller and the buyer ($p + 1/K$ and $p$, respectively). Thus, our estimation procedure GFT-EsT has an asymmetric structure (see the pseudocode, in particular Lines 17 and 20).

First, GFT-EsT draws a sample from a Bernoulli distribution with parameter $(pK+1)/(K+1)$ (Line 15). If the result is 1, it posts price $p$ to the buyer, and the seller receives a price drawn uniformly at random from $[0, p + 1/K]$ (Line 17). Otherwise, if the result is 0, GFT-EsT posts price $p$ to the seller, and the buyer's price is drawn uniformly at random from $[p, 1]$. We denote the final estimate at $t$ by $\widehat{\text{GFT}}_t(p + 1/K, p)$ (Line 20). Overall, our estimator has a small bias, as formalized in the following Lemma.

---

**Algorithm 2:** BLOCK-DECOMPOSITION

---

   **Input:**   • Number of rounds $T$ and number of blocks $N$

             • Set of prices $H_K$

1  Initialize HEDGE over action space $H_K$ and time horizon $N$

2  Initialize random mappings $h_j$ for all $j \in \{0, \dots, N-1\}$

3  $\mathcal{B}_j \leftarrow \{j\frac{T}{N} + 1, \dots, (j+1)\frac{T}{N}\}$ for all $j \in \{0, \dots, N-1\}$

4  **for** $j \in \{0, \dots, N-1\}$ **do**

5      Receive from $\mathcal{A}$ the distribution over pair of prices $x_j$

6      **for** $t \in \mathcal{B}_j$ **do**

7          **if** $t \notin S_j$ **then**

8              Play $(p, q) \sim x_j$ and observe $\mathbb{I}\{s_t \leq p \wedge q \leq b_t\}$

9          **else**

10              Select prices $(p, q)$ such that $h_j(p, q) = t$

11              Compute $\widehat{\text{GFT}}_t(p, q)$ through GFT-EST

12              $\hat{r}_j(p, q) \leftarrow \widehat{\text{GFT}}_t(p, q)$

13      Update $\mathcal{A}$ with reward vector $\hat{r}_j$

---

14  **function** *GFT-EST*

     **Input:** prices $(p + 1/K, p) \in H_K$

15      Sample $Z$ from a Bernoulli with parameter $\frac{pK+1}{K+1}$

16      **if** $Z = 1$ **then**

17          Post price $(\tilde{p}, p)$, with $\tilde{p} \sim U[0, p + 1/K]$

18          $\widehat{\text{GFT}}_t(p + 1/K, p) \leftarrow \mathbb{I}\{s_t \leq \tilde{p}\}\mathbb{I}\{p \leq b_t\}$

19      **else**

20          Post price $(p, \tilde{p})$, with $\tilde{p} \sim U[p, 1]$

21          $\widehat{\text{GFT}}_t(p + 1/K, p) \leftarrow \mathbb{I}\{s_t \leq p\}\mathbb{I}\{\tilde{p} \leq b_t\}$

22      **return** $\widehat{\text{GFT}}_t(p + 1/K, p)$

---

**Lemma 5.2.** *For every $(p + 1/K, p) \in H_K$, the random variable $\widehat{GFT}_t(p + 1/K, p)$ is an $1/K$-biased estimate of $GFT_t(p + 1/K, p)$, i.e.,*

$$\left| GFT_t\left(p + \tfrac{1}{K}, p\right) - \mathbb{E}\left[\widehat{GFT}_t\left(p + \tfrac{1}{K}, p\right)\right]\right| \leq \frac{2}{K}.$$

*Proof.* First, we observe that for $\tilde{p} \sim U[0, p + 1/K]$ (*i.e.,* drawn independently from the uniform distribution over the $U[0, p + 1/K]$ interval) we have that

$$\mathbb{E}[\mathbb{I}\{s_t \leq \tilde{p}\}] = \mathbb{I}\{s_t \leq p + 1/K\}\left(1 - \frac{s_t}{p + 1/K}\right),$$

and for $\tilde{p} \sim U[p, 1]$ we have

$$\mathbb{E}[\mathbb{I}\{\tilde{p} \leq b_t\}] = \mathbb{I}\{p \leq b_t\}\left(\frac{b_t - p}{1 - p}\right).$$

Using these two equations, we can compute the expected value of the random variable returned by GFT-EST. Indeed, by the law of total expectation, we have

18

$$\mathbb{E}\left[\widehat{\mathrm{GFT}}_t\left(p + \frac{1}{K}, p\right)\right]$$

$$= \frac{pK+1}{K+1}\mathbb{I}\{p \le b_t\} \underset{\tilde{p} \sim U[0, p+\frac{1}{K}]}{\mathbb{P}}[s_t \le \tilde{p}] + \frac{1-p}{1 + 1/K}\mathbb{I}\{s_t \le p + 1/K\} \underset{\tilde{p} \sim U[p,1]}{\mathbb{P}}[\tilde{p} \le b_t)]$$

$$= \frac{pK+1}{K+1}\mathbb{I}\{p \le b_t\}\mathbb{I}\{s_t \le p + 1/K\}\left(1 - \frac{s_t}{p + 1/K}\right) + \frac{1-p}{1+1/K}\mathbb{I}\{s_t \le p + 1/K\}\mathbb{I}\{p \le b_t\}\frac{b_t - p}{1-p}$$

$$= \frac{K}{K+1}\left(b_t - s_t + 1/K\right)\mathbb{I}\{s_t \le p + 1/K\}\mathbb{I}\{p \le b_t\}$$

We can thus conclude the proof by observing that:

$$\left|\mathrm{GFT}_t\left(p + \frac{1}{K}, p\right) - \mathbb{E}\left[\widehat{\mathrm{GFT}}_t\left(p + \frac{1}{K}, p\right)\right]\right| = \left|\mathbb{I}\left\{s_t \le p + \frac{1}{K}\right\}\mathbb{I}\{p \le b_t\}\left(b_t - s_t - \frac{b_t - s_t + \frac{1}{K}}{1 + \frac{1}{K}}\right)\right| \le \frac{2}{K},$$

where the last inequality holds since $\left|a - \frac{a+\varepsilon}{1+\varepsilon}\right| \le 2\varepsilon$ for all $a \in [-1, 1]$ and $\varepsilon < 1$. $\qquad\square$

Given the estimation procedure GFT-Est, it is possible to turn any no-regret algorithm for the full-feedback setting into a regret minimizer for the partial feedback setting by the standard block decomposition technique (see, e.g., Chapter 4 of Nisan et al. [2007]). The procedure, which we call Block-Decomposition is described in the pseudocode. We assume to employ Hedge as the full-feedback regret minimizer $\mathcal{A}$.

Block-Decomposition works by subdividing the time horizon $T$ into $N$ blocks $\mathcal{B}_1, \ldots, \mathcal{B}_N$ of equal size and contiguous, that is $\mathcal{B}_j = \{jT/N + 1, \ldots, (j+1)T/N\}$ for any $j \in \{0, 1, \ldots, N-1\}$. In each block we select uniformly at random $K$ time steps (*i.e.,* one for each pair in $H_K$), and we randomly assign each of such time steps to one pair of prices in $H_K$. Formally, for each block $j$, we have a one-to-one map $h_j : \mathcal{B}_j \to H_K$ which is a uniform random map from prices in $H_K$ to rounds in block $\mathcal{B}_j$. We call the image of $h_j$ the *exploration rounds*, and we denote the set of such rounds by $S_j$. For any block $j$, the algorithm builds a vector $\hat{r}_j$ such that the entry $\hat{r}_j(p, q)$ is an estimation of the reward of the pair $(p, q) \in H_K$ in block $\mathcal{B}_j$. To do that, for any block $j$ and pair of prices $(p, q) \in H_K$, we let $\hat{r}_j(p, q) = \widehat{\mathrm{GFT}}_t(p, q)$, where $t = h_j(p, q)$ and $\widehat{\mathrm{GFT}}_t(p, q)$ is computed through the estimation procedure GFT-Est with prices $(p, q)$ (Lines 10 and 11). For any block $j$, exploration rounds in $S_j$ are used to build $\hat{r}_j$. In all the other rounds in $\mathcal{B}_j \setminus S_j$ the algorithm plays according to the strategy $x_j \in \Delta(H_K)$ (Line 8) computed by $\mathcal{A}$ at the beginning of block $j$ (Line 5). At the end of each block $j$, the full-information subroutine $\mathcal{A}$ is updated using $\hat{r}_j$ (Line 13).

Let $\mathrm{GFT}_j(p, q) = \sum_{t \in \mathcal{B}_j} \mathrm{GFT}_t(p, q)/|\mathcal{B}_j|$ be the average GFT over block $\mathcal{B}_j$. Since we choose exploration rounds uniformly at random throughout block $\mathcal{B}_j$ we have that, for any $(p, q) \in H_K$,

$$\left|\mathbb{E}\left[\hat{r}_j(p, q)\right] - \mathrm{GFT}_j(p, q)\right| = \left|\sum_{t \in \mathcal{B}_j} \frac{1}{|\mathcal{B}_j|}\left(\mathbb{E}\left[\widehat{\mathrm{GFT}}_t(p, q)\right] - \mathrm{GFT}_t(p, q)\right)\right|$$

$$\le \sum_{t \in \mathcal{B}_j} \frac{1}{|\mathcal{B}_j|}\left|\mathbb{E}\left[\widehat{\mathrm{GFT}}_t(p, q)\right] - \mathrm{GFT}_t(p, q)\right|$$

$$\le \frac{2}{K},$$

where the last equality follows from Lemma 5.2. This yields the following guarantees on the regret of BLOCK-DECOMPOSITION. Let $x_t$ be the distribution over $H_K$ employed to sample $(p, q)$ at time $t$. At time $t, t \in \mathcal{B}_j$, we have

$$
x_t = \begin{cases} x_j & \text{if } t \notin S_j \text{ (Line 8)} \\ \text{play } (p, q) \text{ s.t. } h_j(p, q) = t & \text{otherwise (Line 11)} \end{cases}
$$

where $x_j$ is the distribution computed by HEDGE for block $j$. The following lemma states precisely the guarantees provided by BLOCK-DECOMPOSITION.

**Lemma 5.3.** BLOCK-DECOMPOSITION *with* $K = T^{1/4}$ *and* $N = T^{1/2}$ *guarantees:*

$$
\sup_{(p,q)\in H_K} \sum_{t=1}^{T} GFT_t(p, q) - \sum_{t=1}^{T} \mathbb{E}_{(p,q)\sim x_t} [GFT_t(p, q)] \leq \frac{5}{2} T^{3/4} \sqrt{\log(T)}.
$$

*Proof.* Let $R_N^H$ be the regret accumulated by HEDGE over $N$ rounds when it observes utilities in $[0, 1]$ and plays over $K$ actions. Each exploration round can cost at most 1 with respect to playing according to $x_j$, and there are $NK$ such rounds. Then, we have that

$$
\sum_{t=1}^{T} \sum_{(p,q)\in H_K} \text{GFT}_t(p, q) x_t(p, q) \geq \sum_{j\in[N]} \sum_{(p,q)\in H_K} \frac{T}{N} \cdot \text{GFT}_j(p, q) x_j(p, q) - NK
$$

$$
\geq \sum_{j\in[N]} \sum_{(p,q)\in H_K} \frac{T}{N} \cdot \left( \mathbb{E}[\hat{r}_j(p, q)] - \frac{2}{K} \right) x_j(p, q) - NK
$$

$$
= \mathbb{E}\left[ \sum_{j\in[N]} \sum_{(p,q)\in H_K} \frac{T}{N} \hat{r}_j(p, q) x_j(p, q) \right] - \frac{2T}{K} - NK
$$

$$
\geq \mathbb{E}\left[ \sup_{(p,q)\in H_K} \sum_{j\in[N]} \frac{T}{N} \hat{r}_j(p, q) \right] - \frac{T}{N} R_N^H - \frac{2T}{K} - NK
$$

$$
\geq \sup_{(p,q)\in H_K} \mathbb{E}\left[ \sum_{j\in[N]} \frac{T}{N} \hat{r}_j(p, q) \right] - \frac{T}{N} R_N^H - \frac{2T}{K} - NK
$$

$$
\geq \sup_{(p,q)\in H_K} \sum_{j\in[N]} \frac{T}{N} \text{GFT}_j(p, q) - \frac{T}{N} R_N^H - \frac{2T}{K} - NK
$$

$$
= \sup_{(p,q)\in H_K} \sum_{t=1}^{T} \text{GFT}_t(p, q) - \frac{T}{N} R_N^H - \frac{2T}{K} - NK.
$$

By rearranging we obtain that

$$
\sup_{(p,q)\in H_K} \sum_{t=1}^{T} \text{GFT}_t(p, q) - \sum_{t=1}^{T} \sum_{(p,q)\in H_K} \text{GFT}_t(p, q) x_t(p, q) \leq \frac{T}{N} R_N^H + \frac{2T}{K} + NK.
$$

20

It is known that $R_N^H \leq 4\sqrt{N \log K}$ (see, *e.g.*, Slivkins [2019]). Then, by setting $K = T^{1/4}$ and $N = T^{1/2}$ we obtain

$$\sup_{(p,q) \in H_K} \sum_{t=1}^{T} \text{GFT}_t(p,q) - \sum_{t=1}^{T} \sum_{(p,q) \in H_K} \text{GFT}_t(p,q) x_t(p,q) \leq T^{3/4} \left( 3 + 4\sqrt{\log(T^{1/4})} \right)$$

$$\leq 16 \cdot T^{3/4} \sqrt{\log(T^{1/4})},$$

where the last inequality holds for all $T \geq 2$. This concludes the proof. $\qquad\square$

### 5.1.3 Putting Everything Together

GFT-Max with the two regret minimizers described in Sections 5.1.1 and 5.1.2 guarantees a $O(T^{3/4})$ bound on the regret.

**Theorem 5.4.** *Consider the repeated bilateral trade problem in the one-bit feedback model. There exists a learning algorithm $\mathcal{A}$ that respects global budget balance and whose regret with respect to the best fixed price in hindsight verifies:*

$$R_T(\mathcal{A}) \leq 1282 \cdot T^{3/4} \log^2 T.$$

*Proof.* The proof follows the same structure of Theorem 4.2. In this case, we set $\beta = T^{3/4}$ and $K = T^{1/4}$, and consider GFT-Max with EXP3.P [Auer et al., 2002] instantiated over set $F_K$ (see Section 5.1.1) as the regret minimizer $\mathcal{A}_P$, and Block-Decomposition instantiated over $H_K$ (see Section 5.1.2) as the regret minimizer $\mathcal{A}_G$.
By construction of Profit-Max, for any stopping time $\tau$ the profit is at least $\beta$, and in rounds $\tau + 1, \ldots, T$ the budget spent is at most $-1/K$. Therefore, the global budget balance condition is satisfied because the cumulative profit at $T$ is at least

$$\beta - (T - \tau)\frac{1}{K} \geq T^{3/4} - T\frac{1}{T^{1/4}} = 0.$$

Now we prove the regret upper bound. For rounds up to $\tau$ we can exploit Lemma 5.1. On the other hand, for any $\tau$, on rounds in $\tau + 1, \ldots, T$ we have

$$\max_{p \in [0,1]} \sum_{t=\tau+1}^{T} \text{GFT}_t(p) \leq \max_{(p,q) \in H_K} \sum_{t=\tau+1}^{T} \text{GFT}_t(p,q) + \frac{T}{K} \tag{11}$$

$$\leq \mathbb{E}_{(p_t,q_t) \sim x_t} \left[ \sum_{t=\tau+1}^{T} \text{GFT}_t(p_t, q_t) \right] + \frac{T}{K} + 5T^{3/4}\sqrt{\log(T^{1/4})}$$

where the first inequality follows from Proposition 3.1, and the second follows by Lemma 5.3 by replacing $T$ with $T - \tau$. Then, assume that the regret bound Lemma 5.1 holds, which happens with probability at least $1 - 1/T$. By summing Equation (11) and Lemma 5.1 we have that, with

21

probability at least $1 - 1/T$,

$$\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \leq \max_{p \in [0,1]} \sum_{t=1}^{\tau} \mathrm{GFT}_t(p) + \max_{p \in [0,1]} \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p)$$

$$\leq 8 \log T(\beta + 1) + \frac{5T}{K} + 256 \log T \sqrt{|F_K| T \log(|F_K| T)} + \max_{p \in [0,1]} \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p)$$

$$\leq \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] + 8 \log T(\beta + 1) + \frac{6T}{K} + 256 \log T \sqrt{|F_K| T \log(|F_K| T)} + 5T^{3/4} \sqrt{\log(T^{1/4})},$$

where the second inequality follows from Lemma 5.1 and the third one from Lemma 5.3.
Then, by substituting $\beta = T^{3/4}$ and $K = T^{1/4}$ we obtain

$$\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) \leq \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] + 8 \log T(T^{3/4} + 1) + 6T^{3/4}$$

$$+ 256 \cdot \log T \sqrt{(2T^{1/4}(\log T + 1))T \log((2T^{1/4}(\log T + 1))T)} + 5T^{3/4} \sqrt{\log(T^{1/4})}$$

$$\leq \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] + 1280 \cdot T^{3/4} \log^2 T.$$

Then, by rearranging, with probability at least $1 - 1/T$ it holds

$$\max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) - \mathbb{E}\left[ \sum_{t=1}^{T} \mathrm{GFT}_t(p_t, q_t) \right] \leq \max_{p \in [0,1]} \sum_{t=1}^{T} \mathrm{GFT}_t(p) - \mathbb{E}\left[ \sum_{t=\tau+1}^{T} \mathrm{GFT}_t(p_t, q_t) \right]$$

$$\leq 1280 \cdot T^{3/4} \log^2 T,$$

where the first inequality follows from the fact that the gain from trade is always non-negative.
Finally, the expected regret is at most

$$R_T(\text{GFT-Max}) \leq \left(1 - \frac{1}{T}\right)\left(1280 \cdot T^{3/4} \log^2 T\right) + \frac{1}{T} \cdot 2T \leq 1282 \cdot T^{3/4} \log^2 T.$$

This concludes the proof. □

## 5.2 $\Omega(T^{5/7})$ Lower Bound with Two-Bit Feedback

In this section, we provide a lower bound for learning the best price against any oblivious adversary, with global budget balance constraints and two-bit feedback. Our construction builds upon the one by Cesa-Bianchi et al. [2023], but exhibits two key differences. First, we are not constrained to use smooth value distributions. This allows us to simplify the construction, avoiding the reduction to online learning with feedback graphs. Second, we only require algorithms to be globally budget balanced (instead of per-round weakly budget balanced); looser budget balance constraints enhance the capabilities of the learning algorithm. All in all, we derive a lower bound that is slightly looser $T^{5/7} \approx T^{0.714}$ compared to the $\Omega(T^{3/4})$. We further elaborate on this comparison at the end of the Section.

**Theorem 5.5.** *Consider the problem of repeated bilateral trade in the two-bit feedback model. Any learning algorithm that satisfies global budget balance suffers regret at least $\Omega(T^{5/7})$.*

The rest of the Section is devoted to the proof of Theorem 5.5; for the missing details, we refer to Appendix A.2. Our lower bound construction is based on $N$ stochastic sequences of valuations. Each one of these sequences is sampled in an i.i.d. way from distributions of valuations with two key properties: (i) they are close with respect to some statistical measure of distance (see Lemma 5.11) and (ii) ensure that any pair of prices that reveals information on the underlying instance is highly suboptimal in terms of GFT (*i.e.*, gathering information is "costly", see Lemma 5.8). We proceed in 5 steps.

**i) Building a set of hard instances.** We start by introducing a set of $N = N(T)$, to be specified later, hard instances of the bilateral trade problem. Our goal is to show that any learning algorithm has regret at least $\Omega(T^{5/7})$ in at least one of the $N$ instances. We define a distribution $\mu_k \in \Delta([0,1]^2)$ of valuations $(s, b)$ over $[0,1]^2$ for each $k \in \{0, \dots, N-1\}$, where we have $N-1$ "perturbed" distributions corresponding to indices $k \in \{1, \dots, N-1\}$, and a "base" distribution corresponding to $k = 0$.

Let $\ell = 1/12$ and let $\Delta = \ell/(N-1)$, and $\delta = \Delta/2$. Then, for any instance $k \in \{0, \dots, N-1\}$, the distributions $\mu_k$ are supported on the same set $\mathcal{W}$ of finitely many valuations. We describe the set $\mathcal{W}$ by partitioning it into six different sets. An illustration of the valuations set can be found in Figure 3a. First, we define the two sets $\mathcal{W}_1$ and $\mathcal{W}_2$ (respectively red and blue in Figure 3a) as follows:

$$\mathcal{W}_1 = \left\{ w_1^i = \left( \tfrac{1-\ell}{2} + i\Delta, 1 - \ell \right) : i = 0, \dots, N-1 \right\}$$

and

$$\mathcal{W}_2 = \left\{ w_2^i = \left( \tfrac{1-l}{2} + i\Delta, 1 - \ell - \rho \right) : i = 0, \dots, N-1 \right\},$$

where $\rho = 1/32$. These valuations are "balanced out" by the $N$ valuations in $\mathcal{W}_3$ (green in Figure 3a):

$$\mathcal{W}_3 = \left\{ w_3^i = \left( 0, \tfrac{1-\ell}{2} - \delta + i\Delta \right) : i = 0, \dots, N-1 \right\}.$$

Moreover, we have a set $\mathcal{W}_4$ of "deficit-generating" valuations (brown in Figure 3a)

$$\mathcal{W}_4 = \left\{ w_4^i = \left( \tfrac{1-\ell}{2} + i\Delta, \tfrac{1-\ell}{2} - \delta + i\Delta \right) : i = 0, \dots, N-1 \right\},$$

and a single valuation belonging to $\mathcal{W}_5$ (orange in Figure 3a)

$$\mathcal{W}_5 = \left\{ \left( 0, \tfrac{1-\ell}{2} \right) \right\}.$$

We conclude by defining the set $\mathcal{W}_6$ (purple in Figure 3a) of the four "extremal" valuations (in practise, they are needed for Lemma 5.11 to hold):

$$\mathcal{W}_6 = \{ (0,0), (0,1), (1,1), (1,0) \}.$$

We assign different probabilities to the valuations in each set $\mathcal{W}_j$ depending on the instance. In particular, for any instance $k \in \{1, \dots, N-1\}$ with distribution $\mu_k$, we have that

$$\mu_k(w_j^i) = \frac{1}{64N^2} = \gamma_1, \quad \forall j \in \{1, 2\}, i \notin \{k, k+1\}, \tag{12}$$

while we perturb by $\varepsilon$ the probability of the following valuations:

$$\mu_k(w_1^k) = \gamma_1 + \varepsilon, \quad \mu_k(w_1^{k+1}) = \gamma_1 - \varepsilon, \quad \mu_k(w_2^k) = \gamma_1 - \varepsilon, \quad \mu_k(w_2^{k+1}) = \gamma_1 + \varepsilon. \tag{13}$$
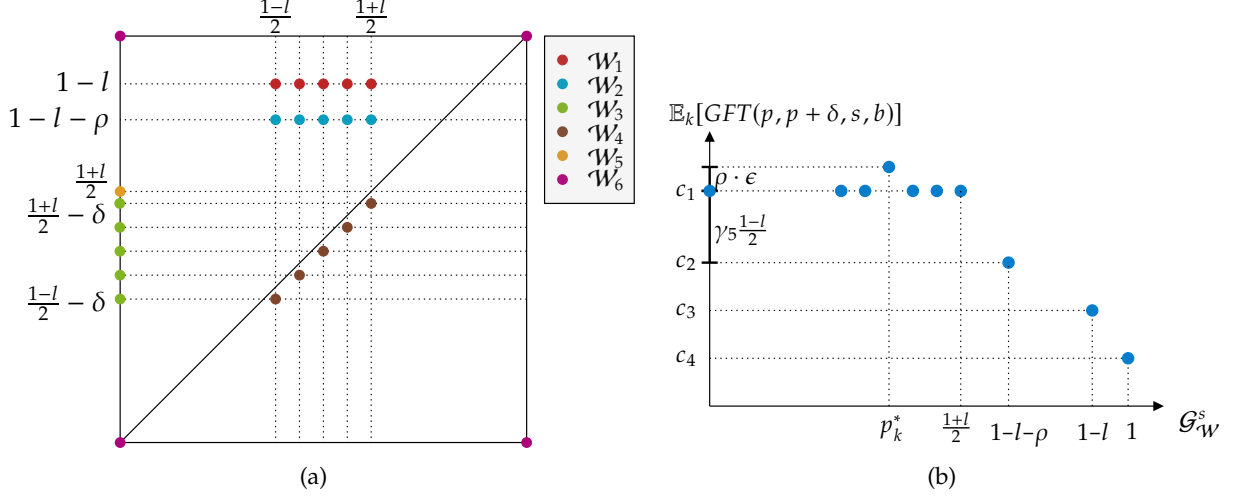
Figure 3: Figure 3a represents the valuations support $\mathcal{W}$ of the instances distributions $\mu_k$, while Figure 3b represents the value of posting the same price to the seller and the buyer in instance $\mu_k$.

Conversely, for the base instance $\mu_0$, we place equal probability $\mu_0(w) = \gamma_1$ on all the valuations $w \in \mathcal{W}_1 \cup \mathcal{W}_2$, and hence all these valuations have the same probability. For each instance $k \in \{0, \ldots, N-1\}$ with distribution $\mu_k$, the probability of valuations $w_3^i$, with $i \in \{0, \ldots, N-1\}$, is set as

$$\mu_k(w_3^i) = \gamma_1 \cdot \frac{1 - \ell - \rho - 2i\Delta}{\frac{1-\ell}{2} - \delta + i\Delta} \in (0, 2\gamma_1).$$

Let $\gamma_3^{\text{tot}} = \sum_{w \in \mathcal{W}_3} \mu_k(w) < 2\gamma_1 N$ be the total probability assigned to valuations in $\mathcal{W}_3$. Moreover, for any instance $k \in \{1, \ldots, N-1\}$ with distribution $\mu_k$, we assign to every point in $\mathcal{W}_4$ probability $\gamma_4 = 4\gamma_1(13N - 14)$, i.e.,

$$\mu_k(w) = 4\gamma_1(13N - 14), \quad \forall w \in \mathcal{W}_4. \tag{14}$$

Then, for any instance $k \in \{1, \ldots, N-1\}$ with distribution $\mu_k$, we assign probability $\gamma_5$ to the single valuation in $\mathcal{W}_5$, i.e., $\mu_k(0, {}^{(1-\ell)}/_2) = \gamma_5 = {}^1/_{64}$. Finally, all the remaining probability is equally divided into the 4 extremal points in $\mathcal{W}_6$, i.e.,

$$\mu_k(w) = \frac{1 - \left(2\gamma_1 N + \gamma_3^{\text{tot}} + 4\gamma_1 N(13N - 14) + \gamma_5\right)}{4} = \gamma_6, \quad \forall w \in \mathcal{W}_6.$$

In Appendix A.2, Lemma A.5 shows that this probabilities are positive and, therefore, $\mu_k$ defines a probability distribution for every $k$.

Now, we define $\mathcal{G}_\mathcal{W}$ as the grid generated by such valuations. Formally:

$$\mathcal{G}_\mathcal{W}^s = \{s \,:\, \exists (s, \cdot) \in \mathcal{W}\}, \mathcal{G}_\mathcal{W}^b = \{b \,:\, \exists (\cdot, b) \in \mathcal{W}\}, \quad \text{and} \quad \mathcal{G}_\mathcal{W} = \left\{(s, b) \,:\, s \in \mathcal{G}_\mathcal{W}^s \text{ and } b \in \mathcal{G}_\mathcal{W}^b\right\}.$$

Thus, $\mathcal{G}_\mathcal{W}^s$ and $\mathcal{G}_\mathcal{W}^b$ represent the projections of $\mathcal{G}_\mathcal{W}$ onto its first (seller) and second (buyer) component, respectively.

**ii) Analysis of the gain from trade.** As a first step, we argue that we can focus on algorithms that play only actions in $\mathcal{G}_\mathcal{W}$, without loss of generality. Consider infact any instance $k \in \{1, \ldots, N\}$ and any algorithm $\mathcal{A}$. Similarly to the proof of Proposition A.3 (more specifically Claim A.4 therein), one can easily prove that there exists an equivalent algorithm $\mathcal{A}'$ (in terms of both feedback, GFT, and profit), that only has distribution supported on the grid $\mathcal{G}_\mathcal{W}$ generated by the valuations $\mathcal{W}$.

Next, for any $p \in \mathcal{G}^s_{\mathcal{W}}$, we characterize the value of posting the pair of prices $(p, p + \delta)$ under distribution $\mu_k$, with $k \in \{0, \ldots, N - 1\}$. Note that posting the pair $(p, p + \delta) \in \mathcal{G}_{\mathcal{W}}$ under any instance $\mu_k$, is equivalent to posting a single price $p \in \mathcal{G}^s_{\mathcal{W}}$ to both the seller and the buyer, with the only difference that $(p, p) \notin \mathcal{G}_{\mathcal{W}}$, while $(p, p + \delta) \in \mathcal{G}_{\mathcal{W}}$. Then, for any $p \in \mathcal{G}^s_{\mathcal{W}}$, we relate the GFT obtained by posting a pair $(p, p + \delta)$ under valuations sampled from $\mu_k$, with $k \in \{1, \ldots, N - 1\}$, and under the base distribution $\mu_0$. For every $k \in \{0, \ldots, N - 1\}$, let $\mathbb{E}_k$ and $\mathbb{P}_k$ denote the expectation and the probability measure under instance $\mu_k$, respectively. Direct calculations shows that, for all $p \in \mathcal{G}^s_{\mathcal{W}}$ and $k \in \{1, \ldots, N - 1\}$, it holds

$$\mathbb{E}_k[\text{GFT}(p, p + \delta, s, b)] = \mathbb{E}_0[\text{GFT}(p, p + \delta, s, b)] + \rho \varepsilon \mathbb{I}\{p = p^*_k\},$$

where $\text{GFT}(p, p + \delta, s, b)$ is simply the gain from trade when the prices posted are $(p, p + \delta)$ and valuations $(s, b)$, and $p^*_k = \frac{1 - \ell}{2} + k\Delta$. Moreover, for all $p \in \mathcal{G}^s_{\mathcal{W}}$ it holds:

$$\mathbb{E}_0[\text{GFT}(p, p + \delta, s, b)] = \begin{cases} c_1 = \gamma_5 \frac{1 + \ell}{2} + \mu_0(0, 1) + \gamma_1 \frac{77}{96} N = & \text{if } p \in [0, \frac{1 + \ell}{2}] \\ c_2 = \mu_0(0, 1) + \gamma_1 \frac{77}{96} N & \text{if } p \in (\frac{1 + \ell}{2}, 1 - \ell - c] \\ c_3 = \mu_0(0, 1) + \gamma_1 \frac{5}{12} N & \text{if } p \in (1 - \ell - c, 1 - \ell] \\ c_4 = \mu_0(0, 1) & \text{if } p \in (1 - \ell, 1] \end{cases}$$

Figure 3b gives a representation of $\mathbb{E}_k[\text{GFT}(p, p + \delta, s, b)]$. From these calculations, we show that in an instance $k \in \{1, \ldots, N - 1\}$ the pair that maximizes the expected gain from trade is $(p^*_k, p^*_k + \delta)$.

**Lemma 5.6.** *For any instance $k \in \{1, \ldots, N - 1\}$, we have that:*

$$\max_{(p,q) \in [0,1]^2,\, p \leq q} \mathbb{E}_k[\text{GFT}(p, q, s, b)] = \mathbb{E}_k[\text{GFT}(p^*_k, p^*_k + \delta, s, b)] = c_1 + \rho \cdot \varepsilon.$$

The previous lemma characterizes the optimal fixed budget balanced price. Then, we show that all the strategies that are *not* budget balanced are dominated. Indeed, one of the main challenges of our reduction is that, in general, a globally budget balanced algorithm could get a larger GFT by temporarily sacrificing some profit and posting prices $(p, q)$ with $q < p$. In the following lemma we show that our instances are built in such a way that these strategies are dominated and thus can be discarded. Intuitively, every tuple of prices $p, q$ that tries to gain higher GFT than the one obtained by playing on the diagonal must win also trades in $\mathcal{W}_4$. Then, since trades in $\mathcal{W}_4$ have negative GFT and happen with sufficiently high probability $\gamma_4$, we have that posting prices $q < p$ is dominated.

**Lemma 5.7.** *For every pair of posted prices $(p, q) \in \mathcal{G}_{\mathcal{W}} \cap \{(p, q) \in [0, 1]^2 \mid p < q\}$, $(p', q') \in \mathcal{G}_{\mathcal{W}} \cap \{(p, q) \in [0, 1]^2 \mid p \geq q\}$, and instance $k \in \{0, \ldots, N - 1\}$, we have that*

$$\mathbb{E}_k[\text{GFT}(p, q, s, b)] \leq \mathbb{E}_k[\text{GFT}(p', q', s, b)] \leq c_1 + \rho \varepsilon \mathbb{I}\{(p', q') = (p^*_k, p^*_k + \delta)\}.$$

We complete this section by showing that also strategies that propose a high price to the buyer are dominated in every instance. In particular, we show that when the algorithm places prices $(p, q)$ with $q > \frac{(1 + \ell)}{2}$, it looses a constant GFT with respect to choosing a smaller $q$. This is because the learner cannot induce the trade $\mathcal{W}_5$ which guarantees expected GFT of $\Theta(\gamma_5)$. Formally,

**Lemma 5.8.** *For any instance $k \in \{0, \ldots, N - 1\}$, price $p \in \mathcal{W}^s_{\mathcal{G}} \cap \left[\frac{1 - \ell}{2}, \frac{1 + \ell}{2}\right]$, and price $q \in \left(\frac{1 + \ell}{2}, 1\right] \cap \mathcal{G}^b_{\mathcal{W}}$ we have that*

$$\mathbb{E}_k[\text{GFT}(p, p + \delta)] \geq \mathbb{E}_k[\text{GFT}(p, q)] + \frac{\gamma_5}{3}.$$

Intuitively, the previous lemma shows that exploring is costly. Indeed, as we show in the following paragraph, the algorithm must post $q \geq (1+\ell)/2$ to gain information on the instance, *i.e.*, on the $k$ that determines the instance.

**iii) Analysis of the feedback.** In the two-bit feedback model, for a valuation $(s, b)$ we have that posting prices $(p, q)$ generates the feedback $(\mathbb{I}\{s \leq p\}, \mathbb{I}\{q \leq b\})$. Now, we show that for any instance $\mu_k$ and any posted prices $(p, q)$, the distribution of the feedback is independent on the instance almost everywhere. Specifically, the feedback distribution depends on the instance $k$ only within a "small" and instance-dependent region of prices. For every instance $k \in \{1, \ldots, N-1\}$, let

$$\mathcal{F}_k = \left[ \tfrac{1-\ell}{2} + (k-1)\Delta, \tfrac{1-\ell}{2} + k\Delta \right) \times (1 - \ell - c, 1 - \ell].$$

It is a simple exercise to see that, for each pair of prices outside the sets $\mathcal{F}_k$, the feedback received by the learner is independent of the specific instance that is generating the valuations (see [Cesa-Bianchi et al., 2023, Claim 2] for a similar result).

**Lemma 5.9.** *For all $(p, q) \in [0, 1]^2 \setminus \bigcup_{k' \in \{1, \ldots, N-1\}} \mathcal{F}_{k'}$ it holds:*

$$\mathbb{P}_k[(\mathbb{I}\{s \leq p\}, \mathbb{I}\{q \leq b\}) = z] = \mathbb{P}_j[(\mathbb{I}\{s \leq p\}, \mathbb{I}\{q \leq b\}) = z], \quad \forall z \in \{0, 1\}^2, \forall j, k \in \{0, \ldots, N-1\}.$$

**iv) Price regions.** The properties uncovered so far naturally partition the square $[0, 1]^2$ into the following three regions:

- *Exploration regions.* We have the $N-1$ regions $\mathcal{F}_k$. These are the regions in which the probability of observing a certain two-bit feedback depends on the instance $\mu_k$ from which the valuations are sampled.

- *Exploitation regions.* We define the regions $\mathcal{E}_k$ for any $k \in \{1, \ldots, N-1\}$ as follows

$$\mathcal{E}_k = \left\{ (p, q) \in [0, 1]^2 \;\middle|\; q \geq p, \; q \leq \tfrac{1+\ell}{2}, \; p \in \left[ \tfrac{1-\ell}{2} + (k-1)\Delta, \tfrac{1-\ell}{2} + k\Delta \right) \right\}.$$

  All these regions are such that the GFT collected by posting $(p, q) \in \mathcal{E}_k$ is close (and smaller than or equal to) to the optimal GFT, *i.e.*, the one obtained by posting $(p_k^*, p_k^* + \delta)$.

- *Dominated regions.* We define $\mathcal{D}$ as the remaining set of possible valuations, that is

$$\mathcal{D} = [0, 1]^2 \setminus (\cup_k (\mathcal{F}_k \cup \mathcal{E}_k)).$$

  It's easy to verify that by posting $(p, q) \in \mathcal{D}$ one obtains a GFT that is at most $c_1$.

Figure 4 shows the partition of the square $[0, 1]^2$ into exploration, exploitation and dominated figures, which are depicted in **red**, **orange** and **green**, respectively. Next, we define

$$\mathcal{N}_k = \sum_{t \in [T]} \mathbb{I}\{(p_t, q_t) \in \mathcal{F}_k\}, \quad \mathcal{M}_k = \sum_{t \in [T]} \mathbb{I}\{(p_t, q_t) \in \mathcal{E}_k\}, \quad O = \sum_{t \in [T]} \mathbb{I}\{(p_t, q_t) \in \mathcal{D}\},$$

which are the number of times an algorithm plays in the exploration, exploitation and dominated regions, respectively. Then, we can upper bound the gain from trade of an algorithm $\mathcal{A}$ considering only the number of times $\mathcal{A}$ plays in each region. In particular, it holds that in any instance $k$:

- **Cost of exploration:** the GFT collected by posting prices in $\mathcal{F}_j$ is at most $c_2$ for all $j$ (Lemma 5.8);

- **Exploitation:** the GFT collected by posting prices in $\mathcal{E}_j$ is at most $c_1 + \rho \cdot \varepsilon \mathbb{I}\{j = k\}$ (Lemma 5.7);
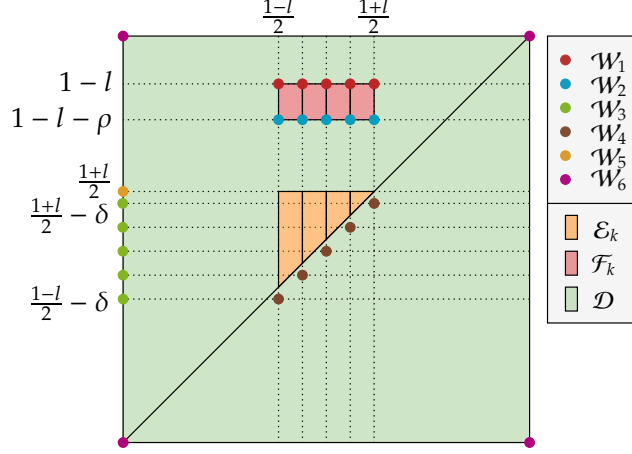
Figure 4: Partition of $[0,1]^2$ in exploration regions $\mathcal{F}_k$, exploitation regions $\mathcal{E}_k$, and dominated regions $\mathcal{D}$.

- **Cost of domination:** the GFT collected by posting prices in $\mathcal{D}$ is at most $c_1$ (Lemma 5.7).

Formally, these observations lead to the following upper bound.

**Lemma 5.10.** *Let* $\{(p_t, q_t)\}_{t \in [T]}$ *be the sequences of prices posted by any algorithm* $\mathcal{A}$. *Then*

$$\sum_{t=1}^{T} \mathbb{E}_k[GFT(p_t, q_t, s, b)] \le \mathbb{E}_k\left[\rho\varepsilon \cdot \mathcal{M}_k + \sum_{k=1}^{N-1}\left(c_1 \mathcal{M}_j + c_2 \mathcal{N}_j + c_1 O\right)\right].$$

**v) Relating the algorithm behavior on different instances.** Now we relate the expected number of exploitation rounds $\mathcal{M}_k$ in different instances $k$. This difference depends on the probability measures $\mathbb{P}_k$ and $\mathbb{P}_0$ through the Pinsker's inequality on a suitably defined multinomial random variable that encodes the four possible feedback observed when playing in the exploration regions $\mathcal{E}_k$.

**Lemma 5.11.** *For all* $k \in \{1, \dots, N-1\}$ *we have that*

$$\mathbb{E}_k[\mathcal{M}_k] - \mathbb{E}_0[\mathcal{M}_k] \le T\varepsilon\sqrt{\frac{2}{\gamma_6}\mathbb{E}_0[\mathcal{N}_k]}.$$

**vi) Lower bounding the regret.**
We define the expected regret of an algorithm under instance $k \in \{0, \dots, N-1\}$ as:

$$R_T^k = \max_{(p,q)\in[0,1]^2, p\ge q} \mathbb{E}_k\left[\sum_{t=1}^{T} \mathrm{GFT}_t(p, q) - \sum_{t=1}^{T} \mathrm{GFT}_t(p_t, q_t)\right].$$

Then, combining all the previous results leads to the following lemma which gives a lower bound in terms of $\varepsilon$, $N$, and $T$.

**Lemma 5.12.** *There is an instance* $k \in \{0, \dots, N-1\}$ *and an absolute constant* $c \in (0,1)$ *such that:*

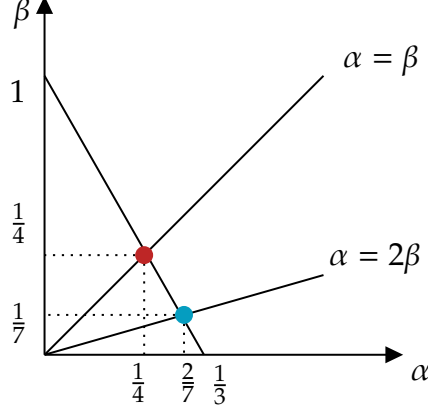$$R_T^k \ge c \cdot \min\left(\tfrac{N}{\varepsilon^2}, \varepsilon T\right).$$

27

Figure 5: Order of $\alpha$ and $\beta$ reachable by Cesa-Bianchi et al. [2023] (**red**) and this work (**blue**).

By using Lemma 5.12 we can readily conclude the proof of Theorem 5.5 as follows. Let $\varepsilon = T^{-\alpha}$ and $N = T^{\beta}$, with $\alpha, \beta > 0$. Now we simply have to optimize over the choice of parameters $\alpha$ and $\beta$. In doing so, we need to take into account the additional constraints necessary to have well-defined instance distributions $\mu_k$. In particular, we have that $\varepsilon \leq \gamma_1$ from Equation (13), and $2N\gamma_1 < 1$ from Equation (12). Moreover, we also need to impose $\gamma_4 N < 1$ by Equation (14). Since $\gamma_4 = 4\gamma_1(14N - 13)$, this also implies that $\gamma_1 < 1/4N(13N-14) < 1/N^2$ for $N > 2$. Therefore, the constraint $\varepsilon < \gamma_1$ implies: $\varepsilon = T^{-\alpha} \leq 1/T^{2\beta} = 1/N^2$ which yields that $\alpha \geq 2\beta$. Note that this dominates the constraint $\varepsilon < 1/N$ (or equivalently written as $\alpha \geq \beta$) that would have been implied by Equation (13) alone.

The lower bound of is maximized when $\alpha$ and $\beta$ are solution of the following program:

$$\max_{\alpha \geq 0} (1 - \alpha)$$
$$\text{s.t.} \quad \alpha \geq 2\beta, \ 1 - \alpha = \beta + 2\alpha$$

which gives as solution the values $\alpha = 2/7$ and $\beta = 1/7$. This implies a $\Omega(T^{5/7})$ lower bound.

**Connection with the $\Omega(T^{3/4})$ lower bound of Cesa-Bianchi et al. [2023].** While our result and the one of Cesa-Bianchi et al. [2023] build on a similar constructions (at least conceptually), we obtain a weaker lower bound. The main reason is that the learner in Cesa-Bianchi et al. [2023] is weak budget balanced, while in our work the learner has only a global budget balance constraint. To preclude this option to the learner, we penalize the GFT of prices in the lower triangle by adding the set of valuations $\mathcal{W}_4$. If $\gamma_4$ is large enough w.r.t. $\gamma_1$, then posting prices in the lower triangle is dominated. In particular, we must choose $\gamma_4 = \Theta(\gamma_1 N)$ as we prove in Lemma 5.7. Once we prove that the lower triangle is dominated, we can conceptually reduce our problem to the one of Cesa-Bianchi et al. [2023]. However, the choice of $\gamma_4 = \Theta(\gamma_1 N)$ imposes the additional constraint $\alpha \geq 2\beta$, which is not needed in the original construction. Hence, they can set $\alpha = \beta = 1/4$, and get a bound of $\Omega(T^{3/4})$. This difference is depicted in Figure 5.

## 6 Best Feasible Distribution of Prices

In this section, we analyse the regret with respect to the best fixed distribution over prices which satisfies global budget balance *on average*. First, we present a negative result that clearly separates

this new benchmark from the best fixed price in hindsight: in Theorem 6.2, we prove that it is impossible to achieve sublinear $(1 + \varepsilon)$-regret with respect to the best feasible distribution, even in the full feedback setting. On the positive side, we show that the two benchmarks are only a multiplicative factor 2 apart (Theorem 6.3). This implies that any learning algorithm that exhibits sublinear regret with respect to the best fixed price in hindsight automatically achieves sublinear 2-regret with respect to the best feasible distribution. Finally, we complement this positive result by proving that this multiplicative gap of 2 is tight (Theorem 6.5).

## 6.1 Linear Lower Bound

The best feasible distribution has a crucial advantage with respect to any budget balanced learner: it has the possibility to "run some deficit" in a preliminary phase of the sequence as it knows it will be possible to extract enough profit to ensure global budget balance in some later stages. For instance, consider a half-sequence where $(s_t, b_t)$ is either $(0, 1/3)$ or $(2/3, 1)$, for $t \leq T/2$. Any learning algorithm has to enforce budget balance at time $T/2$ (to be protected about the possibility that $(s_t, b_t) = 0$ for all future $t$), while the randomized benchmark, which knows the future, may run a deficit and collect more gain from trade by posting the budget unbalanced prices $(2/3, 1/3)$ with some probability. Inspired by this example, we state the following Lemma.

**Lemma 6.1.** *For any algorithm $\mathcal{A}$ that enforces global budget balance, there exists a deterministic sequence of valuations $\mathcal{S}_1$ with the following properties: (i) the expected gain from trade of $\mathcal{A}$ is at most $T/9$; (ii) the valuations $(s_t, b_t)$ are either $(0, 1/3)$ or $(2/3, 1)$ for all $t \leq T/2$; (iii) the valuations $(s_t, b_t)$ are equal to $(0,0)$ for all $t > T/2$.*

*Proof.* Consider the following randomized instance: $(s_t, b_t) = (0, 0)$ for $t > T/2$, while for the other time steps the valuations are either $(0, 1/3)$ or $(2/3, 1)$, independently and uniformly at random. Each realized instance of such randomized sequence clearly satisfies requirements $(ii)$ and $(iii)$. Finally, we show that in expectation (with respect to the randomization of the algorithm and of the instance), the total gain from trade of $\mathcal{A}$ is at most $T/9$. Then the existence of an instance $\mathcal{S}_1$ with the desired properties follows by an averaging argument.
Focus on the first $T/2$ time steps, and let $N_1$ be the random variable that counts the number of time steps in which $\mathcal{A}$ posts prices $q_t \leq 1/3$ and $p_t \geq 2/3$. Moreover, let $N_2 = T/2 - N_1$, and $n_i = \mathbb{E}[N_i]$ for $i \in \{1, 2\}$. By assumption, Algorithm $\mathcal{A}$ is global budget balanced, which this means that

$$-\frac{n_1}{3} + \frac{n_2}{6} \geq 0.$$

The first term of the inequality follows from the fact that every time the learner posts $q_t \leq 1/3$ and $p_t \geq 2/3$, it loses at least $1/3$ revenue. The second term follows from the fact that, by posting other pairs of prices, the learner can extract at most a revenue of $1/6$ (*i.e.,* the trade happens with probability $1/2$, and the learner receives $1/3$ revenue). On the other hand, $n_1$ is directly proportional to the final gain from trade of $\mathcal{A}$, thus the best possible gain from trade is achieved for $n_1 = T/6$ and $n_2 = T/3$, which yields an expected gain from trade of at least

$$\frac{n_1}{3} + \frac{n_2}{6} = \frac{T}{9}.$$

This concludes the proof of the claim. $\qquad\square$

The lemma is crucial in proving the impossibility result in the following Theorem, which holds even under full feedback.

**Theorem 6.2.** *Fix any constant $\alpha \in [1, 36/35)$, and any globally budget balanced learning algorithm $\mathcal{A}$ with full-feedback. Then there exists a sequence of valuations such that*

$$\sum_{t=1}^{T} \mathop{\mathbb{E}}_{(p,q)\sim\gamma^*} GFT_t(p,q) - \alpha \cdot \sum_{t=1}^{T} \mathbb{E}[GFT_t(p_t,q_t)] \geq \tfrac{5}{18} \left(\tfrac{36}{35} - \alpha\right) T,$$

*where distribution $\gamma^*$ is the optimal feasible distribution.*

*Proof.* Fix any $\alpha \in [1, 36/35)$ and any learning algorithm $\mathcal{A}$. Starting from sequence $\mathcal{S}_1$ as in Lemma 6.1, construct a second sequence of valuations $\mathcal{S}_2$ which coincides with $\mathcal{S}_1$ for the first half of the time horizon. In the second half we set $(s_t, b_t) = (\hat{s}, \hat{b})$ for all $t > T/2$, where $(\hat{s}, \hat{b})$ is the most frequent value in the first half of $\mathcal{S}_1$. We compare the total gain from trade collected by $\mathcal{A}$ on $\mathcal{S}_2$ with that of the best fixed distribution of prices over $\mathcal{S}_2$, whose gain from trade we denote OPT. The expected gain from trade of $\mathcal{A}$ on $\mathcal{S}_2$ is at most $T/9$ in the first half (Claim 6.1) and $T/6$ in the second half (as it can extract at most $1/3$ gain from trade in each one of $T/2$ time steps). Therefore, $\mathcal{A}$ extracts at most a total of $5T/18$ expected gain from trade. On the other hand, the best feasible distribution $\gamma^*$ must perform at least as well as the feasible distribution $\gamma$, under which the prices are $(\hat{s}, \hat{b})$ with probability $4/7$, and $(2/3, 1/3)$ with the remaining probability. First, we argue that distribution $\gamma$ is indeed budget-feasible:

$$\sum_{t=1}^{T} \mathop{\mathbb{E}}_{(p,q)\sim\gamma} [\text{PROFIT}_t(p,q)] \geq \tfrac{3}{7}\left(-\tfrac{T}{3}\right) + \tfrac{4}{7}\left(\tfrac{1}{3} \cdot \tfrac{3T}{4}\right) = 0,$$

where we used that prices $(2/3, 1/3)$ are posted with probability $3/7$ and always induces a negative profit of $1/3$ and that, by construction, there are at least $3T/4$ time steps where $(s_t, b_t) = (\hat{s}, \hat{b})$. To conclude the proof, we analyze in a similar way the total gain from trade achieved by $\gamma$:

$$\sum_{t=1}^{T} \mathop{\mathbb{E}}_{(p,q)\sim\gamma} [\text{GFT}_t(p,q)] \geq \tfrac{3}{7}\left(\tfrac{T}{3}\right) + \tfrac{4}{7}\left(\tfrac{1}{3} \cdot \tfrac{3T}{4}\right) = \tfrac{2}{7}T.$$

All in all, we have constructed an instance, $\mathcal{S}_2$ where $\mathcal{A}$ exhibits an expected gain from trade of at most $5T/18$, while OPT is at least $2T/7$. This means that $\mathcal{A}$ suffers at least the following $\alpha$-regret

$$\begin{aligned}
\text{OPT} - \alpha \cdot \sum_{t=1}^{T} \mathbb{E}[\text{GFT}_t(p_t,q_t)] &\geq \sum_{t=1}^{T} \mathop{\mathbb{E}}_{(p,q)\sim\gamma} [\text{GFT}_t(p,q)] - \alpha \cdot \sum_{t=1}^{T} \mathbb{E}[\text{GFT}_t(p_t,q_t)] \\
&\geq \left(\tfrac{2}{7} - \alpha\tfrac{5}{18}\right) T = \tfrac{5}{18}\left(\tfrac{36}{35} - \alpha\right) T. \qquad \square
\end{aligned}$$

## 6.2 Comparison of the Two Benchmarks

Surprisingly, it holds that the performance of the optimal fixed price is to not far from that of optimal global budget balanced distribution.

**Theorem 6.3.** *Denote with $p^*$, resp. $\gamma^*$, the best fixed price, resp. the best feasible distribution. Then, for any sequence of valuations:*

$$\sum_{t=1}^{T} \mathop{\mathbb{E}}_{(p,q)\sim\gamma^*} GFT_t(p,q) \leq 2 \sum_{t=1}^{T} GFT_t(p^*).$$

*Proof.* Fix any sequence of valuations $\mathcal{S}$, and let $p^*$ be as in the statement. By standard analytic arguments, it is possible to show that there exists an optimal feasible distribution $\gamma^*$ whose support is either one or two points (we refer to Proposition A.3 in Appendix A.1 for a formal proof). We prove the result using this $\gamma^*$ and considering two separate cases, according to the cardinality of the support of $\gamma^*$. If the support of $\gamma^*$ consists of only one point $(p, q)$, and since $\gamma^*$ has respect budget feasibility, then it is safe to assume without loss of generality that such point lies above the diagonal, *i.e.*, $p \leq q$ and that the gain from trade achieved by $\gamma^*$ is exactly the same provided by $p^*$.

In the second case, the support of $\gamma^*$ consists of two different points $(p_1, q_1)$ and $(p_2, q_2)$. If both prices lie in the upper left diagonal (i.e., $p_1 \leq q_1$ and $p_2 \leq q_2$), then the total gain from trade is exactly the same as $p^*$, by maximality of $p^*$. If one of the two pair of prices is strongly budget balance, let's say $p_1 = q_1$ and $q_2 < p_1$, then the only possibility (by the budget balance condition) is that these prices never incur in negative profit, so that their gain from trade is once again at most that of $p^*$. All in all, the only meaningful case to study is when $p_1 < q_1$ and $p_2 > q_2$.

Consider then this case, i.e., $p_1 < q_1$ and $p_2 > q_2$, let $\mathcal{T}_0$ be the set of time steps in which the trade is lost by $(p_2, q_2)$, that is $\mathcal{T}_0 = \{t \in [T] \mid s_t > p_2 \text{ or } b_t < q_2\}$. For all other $t \in [T] \setminus \mathcal{T}_0$, every prices $(p_2, q_2)$ make the trade happen. We further partition these time steps as follows:

$$\mathcal{T}_1 = \{t : (s_t, b_t) \in [0, p_2] \times (p_2, 1]\}, \mathcal{T}_2 = \{t : (s_t, b_t) \in [0, q_2) \times [q_2, p_2]\}, \mathcal{T}_3 = \{t : (s_t, b_t) \in [q_2, p_2]^2\}.$$

The sets $\mathcal{T}_0, \ldots, \mathcal{T}_3$ partition the time horizon. Now, for each one of these subset of time steps $\mathcal{T}_i$ it is possible to define two functions over $[0, 1]^2$:

$$f_i(p, q) = \sum_{t \in \mathcal{T}_i} \mathrm{GFT}_t(p, q), \quad g_i(p, q) = \sum_{t \in \mathcal{T}_i} \mathrm{PROFIT}_t(p, q).$$

We adopt the usual convention to omit the second argument if it coincides with the first one. Clearly, the sum of the $f_i$ yields the total GFT, while that of the total $g_i$ the PROFIT. We relate the value of functions $f_0, f_1, f_2$ in $(p_2, q_2)$ with the total gain from trade it collects. The trades in $\mathcal{T}_0$ are lost by $(p_2, q_2)$, so it holds that $f_0(p_2, q_2) = 0$ and $g_0(p_2, q_2) = 0$. For $\mathcal{T}_1$ and $\mathcal{T}_2$ these simple bounds hold:

$$f_1(p_2, q_2) + f_2(p_2, q_2) \leq f_1(p_2) + f_2(q_2) \leq 2 \sum_{t=1}^{T} \mathrm{GFT}_t(p^*). \tag{15}$$

We move our attention to $f_3$, where a more sophisticated argument is needed. As a preliminary step, we prove that the profit extracted by $(p_1, q_1)$ is at most the optimal gain from trade:

$$\sum_{t=1}^{T} \mathrm{PROFIT}_t(p_1, q_1) = (q_1 - p_1) \sum_{t=1}^{T} \mathbb{I}\{s_t \leq p_1\} \mathbb{I}\{q_1 \leq b_t\}$$

$$\leq \sum_{t=1}^{T} (b_t - s_t) \mathbb{I}\{s_t \leq p_1\} \mathbb{I}\{q_1 \leq b_t\}$$

$$\leq \sum_{t=1}^{T} \mathrm{GFT}_t(p_1) \leq \sum_{t=1}^{T} \mathrm{GFT}_t(p^*). \tag{16}$$

Let $\pi_1$, respectively $\pi_2$, be the probability with which $\gamma^*$ draws $(p_1, q_1)$, respectively $(p_2, q_2)$ we have:

$$f_3(p_2, q_2) \leq -g_3(p_2, q_2) \leq -\sum_{t=1}^{T} \mathrm{PROFIT}_t(p_2, q_2) \leq \frac{\pi_1}{\pi_2} \sum_{t=1}^{T} \mathrm{PROFIT}_t(p_1, q_1) \leq \frac{\pi_1}{\pi_2} \sum_{t=1}^{T} \mathrm{GFT}_t(p^*), \tag{17}$$

31

where the first inequality follows by the definition of $\mathcal{T}_3$, the second by the fact that the only negative profit by posting $(p_2, q_2)$ comes from $\mathcal{T}_3$, the third by global budget balance of $\gamma$, and the last one by Equation (16). We finally have all the ingredients to conclude the proof:

$$
\sum_{t=1}^{T} \mathbb{E}_{(p,q)\sim\gamma^*}[\mathrm{GFT}_t(p,q)] = \pi_1 \sum_{t=1}^{T} \mathrm{GFT}_t(p_1,q_1) + \pi_2 \sum_{i=0,\dots,3} f_i(p_2,q_2) \qquad \text{(linearity of exp.)}
$$

$$
\leq \pi_1 \sum_{t=1}^{T} \mathrm{GFT}_t(p_1,q_1) + (\pi_1 + 2\pi_2) \sum_{t=1}^{T} \mathrm{GFT}_t(p^*) \qquad \text{(by Eq. 15 and 17)}
$$

$$
\leq 2 \sum_{t=1}^{T} \mathrm{GFT}_t(p^*),
$$

where the last inequality follows by optimality of $p^*$ with respect to the budget balanced prices $(p_1, q_1)$ and using that $\pi_1 + \pi_2 = 1$. $\qquad\square$

As a corollary, we have that any algorithm that achieves sublinear regret with respect to the best fixed price also guarantees sublinear 2-regret with respect to the best feasible prices distribution.

**Corollary 6.4.** *Let $\mathcal{A}$ be a learning algorithm for the repeated bilateral trade problem which guarantees an upper bound of $f(T)$ on the regret with respect to the best fixed price in hindsight. Then, the 2-regret of $\mathcal{A}$ with respect to the best budget feasible distribution over prices is at most $f(T)$.*

Surprisingly, the factor 2 between the two benchmarks is optimal. This implies that the analysis of the performance of the algorithms in Corollary 6.4 is essentially tight.

**Theorem 6.5.** *For any $\varepsilon > 0$, there exists a sequence of valuations such that*

$$
\sum_{t=1}^{T} \mathbb{E}_{(p,q)\sim\gamma^*} \mathrm{GFT}_t(p,q) \geq (2 - \varepsilon) \sum_{t=1}^{T} \mathrm{GFT}_t(p^*),
$$

*where $p^*$ and $\gamma^*$ are the best fixed price and global budget balanced distribution, respectively.*

*Proof.* Fix any $\varepsilon > 0$, and let $\delta$ be a positive number we set later. Consider the sequence where $(s_t, b_t) = (0, 1/2 - \delta)$ if $t$ is odd, and $(s_t, b_t) = (1/2 + \delta, 1)$ otherwise. Any fixed price can make at most half of the trades happen, with a total gain from trade of at most $T/2\,(1/2 - \delta)$. Consider now the distribution over prices $\gamma$ selecting $(p_1, q_1) = (1/2 + \delta, 1/2 - \delta)$ with probability $\alpha = (1-2\delta)/(1+6\delta)$, and $(p_2, q_2) = (0, 1/2 - \delta)$ otherwise. We conclude the proof by arguing that $\gamma$ satisfies the budget balance constraints, and attains total gain from trade that is roughly twice that of $p^*$. First, we show that $\gamma$ is global budget balanced. We have

$$
\sum_{t\in[T]} \mathbb{E}_{(p,q)\sim\gamma}[\mathrm{PROFIT}_t(p,q)] = \alpha \sum_{t\in[T]} \mathrm{PROFIT}_t(p_1,q_1) + (1-\alpha) \sum_{t\in[T]} \mathrm{PROFIT}_t(p_2,q_2)
$$

$$
\geq \alpha\,(-2\delta T) + (1-\alpha)\tfrac{T}{2}(\tfrac{1}{2} - \delta) = 0,
$$

where in the last equality we use the definition of $\alpha$. We move our attention to the gain from trade:

$$
\sum_{t=1}^{T} \mathbb{E}_{(p,q)\sim\gamma}[\mathrm{GFT}_t(p,q)] = \alpha \sum_{t=1}^{T} \mathrm{GFT}_t(p_1,q_1) + (1-\alpha) \sum_{t=1}^{T} \mathrm{GFT}_t(p_2,q_2)
$$

$$
\geq \alpha T \left(\tfrac{1}{2} - \delta\right) + (1-\alpha)\tfrac{T}{2} \left(\tfrac{1}{2} - \delta\right)
$$

$$
\geq (1+\alpha) \sum_{t=1}^{T} \mathrm{GFT}_t(p^*).
$$

Plugging in the last formula the definition of $\alpha$ and setting $\delta = \varepsilon/8$ yields the desired result. □

# 7 Final Remarks and Open Problems

In this paper we introduce the notion of global budget balance in the repeated bilateral trade problem. With this notion, we show for the first time that it is possible to achieve sublinear regret with respect to the best fixed price in hindsight, without relying on any additional assumption. In the full feedback model we prove that the minimax regret rate of the learning problem is $\tilde{\Theta}(\sqrt{T})$, while in the partial feedback models, we provide an upper bound on the regret of order $\tilde{O}(T^{3/4})$, which is complemented with a $\Omega(T^{5/7})$ lower bound. Our regret results proves a clear separation between the two feedback models, but leave an open gap between the $T^{5/7}$ and $T^{3/4}$ rates in partial feedback.

Inspired by Bandits with Knapsack, we formulated a new benchmark: the best feasible distribution over prices. Against this harder benchmark we prove that it is possible to achieve sublinear 2-regret, while no algorithm can achieve sublinear $(1 + \varepsilon_0)$-regret. We leave as an open question the characterization of the optimal competitive ratio $\alpha \in [1 + \varepsilon_0, 2]$ obtainable against this benchmark.

# 8 References

Shipra Agrawal and Nikhil R. Devanur. Bandits with global convex constraints and objective. *Oper. Res.*, 67(5):1486–1502, 2019. doi: 10.1287/opre.2019.1840.

Noga Alon, Nicolò Cesa-Bianchi, Claudio Gentile, Shie Mannor, Yishay Mansour, and Ohad Shamir. Nonstochastic multi-armed bandits with graph-structured feedback. *SIAM J. Comput.*, 46(6): 1785–1826, 2017. doi: 10.1137/140989455.

Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002. doi: 10.1137/S0097539701398375.

Yossi Azar, Amos Fiat, and Federico Fusco. An *alpha*-regret analysis of adversarial bilateral trade. In *NeurIPS*, 2022.

Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. Bandits with knapsacks. *Journal of the ACM*, 65(3):1–55, 2018. doi: 10.1145/3164539.

Santiago R. Balseiro and Yonatan Gur. Learning in repeated auctions with budgets: Regret minimization and equilibrium. *Manag. Sci.*, 65(9):3952–3968, 2019. doi: 10.1287/mnsc.2018.3174.

Gábor Bartók, Dean P. Foster, Dávid Pál, Alexander Rakhlin, and Csaba Szepesvári. Partial monitoring - classification, regret bounds, and algorithms. *Math. Oper. Res.*, 39(4):967–997, 2014. doi: 10.1287/moor.2014.0663.

Martino Bernasconi, Matteo Castiglioni, Andrea Celli, Alberto Marchesi, Francesco Trovò, and Nicola Gatti. Optimal rates and efficient algorithms for online bayesian persuasion. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 2164–2183. PMLR, 2023.

Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. Bandits with replenishable knapsacks: the best of both worlds. In *ICLR*, 2024a.

Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. No-regret learning in bilateral trade via global budget balance. In *STOC*. ACM, 2024b.

Liad Blumrosen and Shahar Dobzinski. Reallocation mechanisms. In *EC*, page 617. ACM, 2014. doi: 10.1145/2600057.2602843.

Liad Blumrosen and Yehonatan Mizrahi. Approximating gains-from-trade in bilateral trading. In *WINE*, volume 10123 of *Lecture Notes in Computer Science*, pages 400–413. Springer, 2016. doi: 10.1007/978-3-662-54110-4_28.

Nataša Bolić, Tommaso Cesari, and Roberto Colomboni. An online learning theory of brokerage. *The 23rd International Conference on Autonomous Agents and Multi-Agent Systems*, 2024.

Johannes Brustle, Yang Cai, Fa Wu, and Mingfei Zhao. Approximating gains from trade in two-sided markets via simple mechanisms. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, pages 589–590, 2017. doi: 10.1145/3033274.3085148.

Matteo Castiglioni, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Online bayesian persuasion. In *NeurIPS*, 2020.

Matteo Castiglioni, Andrea Celli, and Christian Kroer. Online learning with knapsacks: the best of both worlds. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 2767–2783. PMLR, 2022.

Matteo Castiglioni, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Regret minimization in online bayesian persuasion: Handling adversarial receiver's types under full and partial feedback models. *Artif. Intell.*, 314:103821, 2023.

Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006. doi: 10.1017/CBO9780511546921.

Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Regret minimization under partial monitoring. *Math. Oper. Res.*, 31(3):562–580, 2006. doi: moor.1060.0206.

Nicolò Cesa-Bianchi, Claudio Gentile, and Yishay Mansour. Regret minimization for reserve prices in second-price auctions. *IEEE Trans. Inf. Theory*, 61(1):549–564, 2015. doi: 10.1109/TIT.2014.2365772.

Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. A regret analysis of bilateral trade. In *EC*, pages 289–309. ACM, 2021. doi: 10.1145/3465456.3467645.

Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. Repeated bilateral trade against a smoothed adversary. In *COLT*, volume 195 of *Proceedings of Machine Learning Research*, pages 1095–1130. PMLR, 2023.

Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. The role of transparency in repeated first-price auctions with unknown valuations. In *STOC*. ACM, 2024.

Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colomboni, Federico Fusco, and Stefano Leonardi. Bilateral trade: A regret minimization perspective. *Mathematics of Operations Research*, 49(1): 171–203, 2024. doi: 10.1287/moor.2023.1351.

Constantinos Daskalakis and Vasilis Syrgkanis. Learning in auctions: Regret is hard, envy is easy. *Games Econ. Behav.*, 134:308–343, 2022.

Yuan Deng, Jieming Mao, Balasubramanian Sivan, and Kangning Wang. Approximately efficient bilateral trade. In *STOC*, pages 718–721. ACM, 2022. doi: 10.1145/3519935.3520054.

Paul Duetting, Guru Guruganesh, Jon Schneider, and Joshua Ruizhi Wang. Optimal no-regret learning for one-sided lipschitz functions. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 8836–8850. PMLR, 2023.

Paul Dütting, Federico Fusco, Philip Lazos, Stefano Leonardi, and Rebecca Reiffenhäuser. Efficient two-sided markets with limited information. In *STOC*, pages 1452–1465. ACM, 2021. doi: 10.1145/3406325.3451076.

Yumou Fei. Improved approximation to first-best gains-from-trade. In *International Conference on Web and Internet Economics*, pages 204–218. Springer, 2022. doi: 10.1007/978-3-031-22832-2_12.

Michal Feldman, Tomer Koren, Roi Livni, Yishay Mansour, and Aviv Zohar. Online pricing with strategic and patient buyers. *Advances in Neural Information Processing Systems*, 29, 2016.

Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of convex analysis*. Springer Science & Business Media, 2004.

Chien-Ju Ho, Aleksandrs Slivkins, and Jennifer Wortman Vaughan. Adaptive contract design for crowdsourcing markets: Bandit algorithms for repeated principal-agent problems. *J. Artif. Intell. Res.*, 55:317–359, 2016.

Nicole Immorlica, Karthik Sankararaman, Robert Schapire, and Aleksandrs Slivkins. Adversarial bandits with knapsacks. *Journal of the ACM*, 69(6):1–47, 2022. doi: 10.1145/3557045.

Rodolphe Jenatton, Jim C. Huang, and Cédric Archambeau. Adaptive algorithms for online convex optimization with long-term constraints. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 402–411. JMLR.org, 2016.

Sham M. Kakade, Adam Tauman Kalai, and Katrina Ligett. Playing games with approximation algorithms. *SIAM J. Comput.*, 39(3):1088–1106, 2009. doi: 10.1137/070701704.

Zi Yang Kang, Francisco Pernice, and Jan Vondrák. Fixed-price approximations in bilateral trade. In *SODA*, pages 2964–2985. SIAM, 2022. doi: 10.1137/1.9781611977073.115.

Thomas Kesselheim and Sahil Singla. Online learning with vector costs and bandits with knapsacks. In *COLT*, volume 125 of *Proceedings of Machine Learning Research*, pages 2286–2305. PMLR, 2020.

Robert D. Kleinberg and Frank Thomson Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *FOCS*, pages 594–605. IEEE Computer Society, 2003. doi: 10.1109/SFCS.2003.1238232.

Raunak Kumar and Robert Kleinberg. Non-monotonic resource utilization in the bandits with knapsacks problem. In *NeurIPS*, 2022.

Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020. doi: 10.1017/9781108571401.

Thodoris Lykouris, Vasilis Syrgkanis, and Éva Tardos. Learning and efficiency in games with dynamic population. In *SODA*, pages 120–129. SIAM, 2016. doi: 10.1137/1.9781611977554.ch17.

Mehrdad Mahdavi, Rong Jin, and Tianbao Yang. Trading regret for efficiency: online convex optimization with long term constraints. *J. Mach. Learn. Res.*, 13:2503–2528, 2012.

Shie Mannor, John N. Tsitsiklis, and Jia Yuan Yu. Online learning with sample path constraints. *J. Mach. Learn. Res.*, 10:569–590, 2009.

R Preston McAfee. The gains from trade under fixed price mechanisms. *Applied economics research bulletin*, 1(1):1–10, 2008.

Jamie Morgenstern and Tim Roughgarden. On the pseudo-dimension of nearly optimal auctions. In *NIPS*, pages 136–144, 2015.

Roger B Myerson and Mark A Satterthwaite. Efficient mechanisms for bilateral trading. *Journal of economic theory*, 29(2):265–281, 1983.

Thomas Nedelec, Clément Calauzènes, Noureddine El Karoui, and Vianney Perchet. Learning in repeated auctions. *Found. Trends Mach. Learn.*, 15(3):176–334, 2022. doi: 10.1561/2200000077.

Noam Nisan, Tim Roughgarden, Éva Tardos, and Vijay V. Vazirani, editors. *Algorithmic Game Theory*. Cambridge University Press, 2007. doi: 10.1017/CBO9780511800481.

José Luis Palacios. On the simple symmetric random walk and its maximal function. *The American Statistician*, 62(2):138–140, 2008. doi: 10.1198/000313008X304846.

Alexander Schrijver. *Combinatorial optimization: polyhedra and efficiency*. Springer Science & Business Media, 2003.

Aleksandrs Slivkins. Introduction to multi-armed bandits. *Found. Trends Mach. Learn.*, 12(1-2):1–286, 2019. doi: 10.1561/2200000068.

Aleksandrs Slivkins, Karthik Abinav Sankararaman, and Dylan J. Foster. Contextual bandits with packing and covering constraints: A modular lagrangian approach via regression. In *COLT*, volume 195 of *Proceedings of Machine Learning Research*, pages 4633–4656. PMLR, 2023.

Wen Sun, Debadeepta Dey, and Ashish Kapoor. Safety-aware algorithms for adversarial contextual bandit. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 3280–3288. PMLR, 2017.

William Vickrey. Counterspeculation, auctions, and competitive sealed tenders. *The Journal of finance*, 16(1):8–37, 1961.

Jonathan Weed, Vianney Perchet, and Philippe Rigollet. Online learning in repeated auctions. In *COLT*, volume 49 of *JMLR Workshop and Conference Proceedings*, pages 1562–1583. JMLR.org, 2016.

Andrew Chi-Chih Yao. Probabilistic computations: Toward a unified measure of complexity (extended abstract). In *FOCS*, pages 222–227. IEEE Computer Society, 1977.

Hao Yu, Michael J. Neely, and Xiaohan Wei. Online convex optimization with stochastic constraints. In *NIPS*, pages 1428–1438, 2017.

Banghua Zhu, Stephen Bates, Zhuoran Yang, Yixin Wang, Jiantao Jiao, and Michael I. Jordan. The sample complexity of online contract design. In *EC*, page 1188. ACM, 2023. doi: 10.1145/3580507. 3597673.

You Zu, Krishnamurthy Iyer, and Haifeng Xu. Learning to persuade on the fly: Robustness against ignorance. In *EC*, pages 927–928. ACM, 2021. doi: 10.1145/3465456.3467593.

# A  Appendix

## A.1  Well Posedness of the Two Benchmarks

**Proposition A.1.** *For any sequence of valuations $\mathcal{S}$ there exists a price $p^* \in [0,1]$ such that:*

$$R_T(\mathcal{A}, \mathcal{S}) = \sum_{t=1}^{T} GFT_t(p^*) - \mathbb{E}\left[\sum_{t=1}^{T} GFT_t(p_t, q_t)\right].$$

*Proof.* Denote the cumulative gain from trade of a pair of price $(p, q)$ as follows:

$$f(p, q) = \sum_{t=1}^{T} GFT_t(p, q). \tag{18}$$

This function is upper semi-continuous on the upper left triangle $\{(p, q) \in [0,1]^2 \mid p \leq q\}$ (see Claim A.2), thus the sup in the definition is indeed a max. For the remaining part of the statement, let $(\hat{p}, \hat{q})$ be any pair of prices in the arg max of Equation (1). It is easy to see that any price $p^* \in [\hat{p}, \hat{q}]$ achieves the same total gain from trade, while trivially respecting the budget balance constraint.  □

**Claim A.2.** *The function $f$ defined in Eq. (18) is upper semi-continuous on $\{(p, q) \in [0,1]^2 \mid p \leq q\}$.*

*Proof of Claim A.2.* The function $f$ is the sum of $T$ terms of the following form:

$$GFT_t(p, q) = \mathbb{I}\{s_t \leq p\}\mathbb{I}\{q \leq b_t\}(b_t - s_t).$$

Moreover, for pairs in $\{(p, q) \in [0,1]^2 \mid p \leq q\}$ the gain from trade is non-zero only for steps $t \in [T]$ such that $s_t \leq b_t$. This implies that $f$ is the sum of at most $T$ step-functions that are upper semi-continuous.  □

**Proposition A.3.** *The definition of the best fixed distribution is well-posed. Moreover, there always exists a feasible distribution $\gamma^*$ with support at most two that attains the* sup.

*Proof.* Fix any sequence of valuations $\{(s_t, b_t)\}_{t=1}^{T}$, we introduce two auxiliary functions:

$$f(p, q) = \sum_{t=1}^{T} GFT_t(p, q) \quad \text{and} \quad g(p, q) = \sum_{t=1}^{T} \text{PROFIT}_t(p, q).$$

We can rewrite the program in Definition 2.1 as:

$$\sup_{\gamma \in \Delta([0,1]^2)} \mathbb{E}_{(p,q) \sim \gamma} [f(p, q)]$$

$$s.t. \quad \mathbb{E}_{(p,q) \sim \gamma} [g(p, q)] \geq 0.$$

As a first step, we show that the support of $\gamma$ can be restricted to a discrete grid $G$. To simplify the exposition, we sort the sets of valuations $\{0, 1, s_1, \ldots, s_T\}$ and $\{0, 1, b_1, \ldots, b_T\}$ in increasing order. Formally, we define the set $\{s^0 = 0, s^1, \ldots, s^T, s^{T+1} = 1\}$, where $s^i \leq s^{i+1}$ for each $i$, and $\{s^i\}_{i=0}^{T+1} = \{s_t\}_{t=1}^{T} \cup \{0, 1\}$. Similarly, we define the set $\{b^0 = 0, b^1, \ldots, b^T, b^{T+1} = 1\}$, where $b^i \leq b^{i+1}$ for each $i$, and $\{b^i\}_{i=0}^{T+1} = \{b_t\}_{b=1}^{T} \cup \{0, 1\}$.[3] The grid $G$ contains all the points of the form $(s^i, b^j)$ with $i, j \in \{0, 1, \ldots, T + 1\}$.

---

[3]For the sake of clarity, we assume that $s_i \neq s_j$ and $s_i, s_j \notin \{0, 1\}$, and $b_i \neq b_j$ and $b_i, b_j \notin \{0, 1\}$ for each $i, j \in [T]$. It is easy to extend our results to the general setting.

Now, we assign any point in $[0,1]^2$ to a point in the grid $G$. In particular, we define the map $\pi_G$ that associates each $(p,q) \in [0,1]^2$ to the upper-most and left-most point in its subset of the partition. Formally, $\pi_G : (p,q) \mapsto (s^i, b^j)$, where $i$ is the greatest index such that $s^i \leq p$, and $j$ is the smallest index such that $b^j \geq q$.

From Claim A.4 we have that in Program (2) we can restrict our attention to distributions $\gamma$ that are supported on $G$ and thus are discrete. Hence, we can rewrite Program (2) as the following linear program:

$$\max_{(x_1,x_2)\in C} x_1 \text{ s.t. } x_2 \geq 0,$$

where $C$ is the convex hull of $\mathcal{X} = \{(f(p,q), g(p,q)) \mid (p,q) \in G\}$. Consider any optimal solution $(x_1^\star, x_2^\star)$ of such linear program. Since $(x_1^\star, x_2^\star)$ belongs to $C$, which is a convex hull of a finite set of points, by Caratheodory's Theorem ( see *e.g.*, Theorem 5.1 of Schrijver [2003]) it can be expressed as a convex combination of 3 points in $\mathcal{X}$.

As a direct implication of first-order optimality conditions (*i.e.*, the gradient $(1,0)$ has to belong to the normal cone of $C$ at $(x_1^\star, x_2^\star)$) we have that $(x_1^\star, x_2^\star)$ must be on the boundary $\partial C$ of $C$. This also yields the existence of an hyperplane supporting $C$ at $(x_1^\star, x_2^\star)$ [see, *e.g.*, Lemma 4.2.1 of Hiriart-Urruty and Lemaréchal, 2004]. Since $C$ is entirely contained in one of the halfspaces defined by the supporting hyperplane, and since $(x_1^\star, x_2^\star) \in \partial C$, it must be the case that either $(x_1^\star, x_2^\star) \in \mathcal{X}$, or we can write the optimal point as a convex combination two points belonging to $\mathcal{X}$ (*i.e.*, the two points defining the face of the polytope containing $(x_1^\star, x_2^\star)$). Call $(p_1, q_1)$ and $(p_2, q_2)$ the two points that are in the preimage of the two points generating $(x_1^\star, x_2^\star)$ according to $f$ and $g$, respectively. We showed that there exists an optimal solution whose support consists of the two (possibly coinciding) points $(p_1, q_1)$ and $(p_2, q_2)$. This concludes the proof. □

**Claim A.4.** *Let $\gamma \in \Delta([0,1]^2)$. There exists a distribution $\gamma_G \in \Delta([0,1]^2)$ with the following three properties:*

- $\mathbb{E}_{(p,q)\sim\gamma} [f(p,q)] = \mathbb{E}_{(p,q)\sim\gamma_G} [f(p,q)];$
- $\mathbb{E}_{(p,q)\sim\gamma} [g(p,q)] \geq \mathbb{E}_{(p,q)\sim\gamma_G} [g(p,q)];$
- $\text{supp}(\gamma_G) \subseteq G.$

*Proof.* We define the distribution $\gamma_G$ on $G$ by assigning to each point of the grid $(s^i, b^j)$, with $i \in \{1,\dots,T\}$ and $j \in \{0,\dots,T-1\}$, the probability mass which $\gamma$ assigns to points in the cell of the grid $\{(p,q) \in [0,1]^2 : s^{i-1} \leq i \leq s^i, b^j \leq q \leq b^{j+1}\}$. Formally, the distribution $\gamma_G$ is such that

$$\mathbb{P}_{(p,q)\sim\gamma_G} \left( (p,q) = (s^i, b^j) \right) = \mathbb{P}_{(p,q)\sim\gamma} \left( \pi_G(p,q) = (s^i, b^j) \right).$$

The new distribution $\gamma_G$ is clearly supported on $G$ and thus verifies the third point of the claim. We now prove the remaining two points. First, by construction, the expected gain from trade is not affected by the change in distribution. Indeed, for each $t$, $p$, and $q$, it holds

$$\text{GFT}_t(p,q) = \text{GFT}_t(\pi_G(p,q))$$

since $\mathbb{I}\{s_t \leq p\}\mathbb{I}\{q \leq b_t\} = 1$ if and only if $\mathbb{I}\{s_t \leq p'\}\mathbb{I}\{q' \leq b_t\} = 1$, where $(p',q') := \pi_G(p,q)$. We conclude the proof by showing that the profit does not decrease. It is sufficient to prove that for each $t$, $p$, and $q$, it holds

$$\text{PROFIT}_t(p,q) \geq \text{PROFIT}_t(\pi_G(p,q)).$$

Since $(q-p) \leq (q'-p')$, $\pi_G(p,q)$ and $(p,q)$ make the same trades happen. Then, $\pi_G(p,q)$ extracts at least the same profit of the pair $(p,q)$. This concludes the proof. □

## A.2 Missing Proofs from Section 5.2

**Lemma A.5.** *The distributions $\mu_k$ are well defined for all $k \in \{0, \ldots, N-1\}$.*

*Proof.* Since for all $w \in \mathcal{W}_1 \cup \mathcal{W}_2 \cup \mathcal{W}_3 \cup \mathcal{W}_4 \cup \mathcal{W}_5$ the weights $\mu_k(w)$ are positive for all $k \in \{0, \ldots, N-1\}$ we just need to prove that $\mu_k(w)$ is positive for all $w \in \mathcal{W}_6$.
Then, using the upper bound on $\gamma_3^{\text{tot}} < 2\gamma_1 N \le 1/32N$ we get that:

$$\gamma_6 \ge \frac{1}{4}\left(1 - \left(\frac{1}{32N} + \frac{1}{32N} + \frac{13N-14}{16N} + \frac{1}{64}\right)\right)$$

$$= \frac{1}{4}\left(1 - \left(\frac{13(N-1)}{16N} + \frac{1}{64}\right)\right) \ge \frac{1}{32}.$$

This, together with the fact that $\gamma_6 \le 1/4$, proves that all the probabilities in the instances' distributions $\mu_k$ are well defined. $\qquad\square$

**Lemma 5.7.** *For every pair of posted prices $(p,q) \in \mathcal{G}_{\mathcal{W}} \cap \{(p,q) \in [0,1]^2 \mid p < q\}$, $(p',q') \in \mathcal{G}_{\mathcal{W}} \cap \{(p,q) \in [0,1]^2 \mid p \ge q\}$, and instance $k \in \{0, \ldots, N-1\}$, we have that*

$$\mathbb{E}_k[GFT(p,q,s,b)] \le \mathbb{E}_k[GFT(p',q',s,b)] \le c_1 + \rho\,\varepsilon\,\mathbb{I}\{(p',q') = (p_k^*, p_k^* + \delta)\}.$$

*Proof.* Consider any point $(p^1, q^1) = \left(\frac{1-\ell}{2} + i\Delta, \frac{1-\ell}{2} + \delta + i\Delta\right)$. For any $i \in \{0, \ldots, N-1\}$ we define $(p^1, q_j^{1'}) = \left(\frac{1-\ell}{2} + i\Delta, \frac{1-\ell}{2} - \delta + (i-j)\Delta\right)$ for each $j \in \{0, \ldots, i\}$. Simple calculations show that:

$$\mathbb{E}_k[GFT(p^1, q^1, s, b) - GFT(p, q_j', s, b)] = \sum_{\iota=j}^{i-1}\left[\mu_k(w_3^\iota)\left(\frac{1-\ell}{2} - \delta + \iota\Delta\right) + \gamma_4\delta\right]$$

$$= \sum_{\iota=j}^{i-1}\left[\gamma_1 \cdot \frac{1-\ell-\rho-2\iota\Delta}{\frac{1-\ell}{2} - \delta + \iota\Delta}\left(\frac{1-\ell}{2} - \delta + \iota\Delta\right) + 4\gamma_1(13N-14)\delta\right]$$

$$= \sum_{\iota=j}^{i-1}\left[\gamma_1(1-\ell-\rho-2\iota\Delta + 4\delta(13N-14))\right],$$

which has to hold for all $i \in \{0, \ldots, N-1\}$. The worst case is when $\iota = N-1$. This yields

$$\gamma_1\left(1 - \ell - \rho - 2\ell + 2\frac{\ell}{N-1}(13N-14)\right) > 0, \quad \forall N > 1$$

and thus for all $j$ we have that:

$$\mathbb{E}_k[GFT(p^1, q^1, s, b) - GFT(p^1, q_j^{1'}, s, b)] > 0.$$

The proof is concluded by noting that, for any $(p,q) \in \mathcal{G}_{\mathcal{W}} \cap \{(p,q) \in [0,1]^2 \mid p < q\}$ in the lower triangle, the gain from trade is upper bounded by that of some $(p^1, q_j^{1'})$, that is

$$\mathbb{E}_k[GFT(p,q,s,b)] \le \mathbb{E}_k[GFT(p^1, q_j^{1'}, s, b)].$$

On the other hand, for any $(p',q') \in \mathcal{G}_{\mathcal{W}} \cap \{(p,q) \in [0,1]^2 \mid p \ge q\}$ in the upper triangle, the gain from trade is lower bounded by that of a point $(p^1, q^1)$, that is

$$\mathbb{E}_k[GFT(p^1, q^1, s, b)] \le \mathbb{E}_k[GFT(p', q', s, b)].$$

This concludes the proof of the lemma. $\qquad\square$

**Lemma 5.11.** *For all $k \in \{1, \ldots, N-1\}$ we have that*

$$\mathbb{E}_k[\mathcal{M}_k] - \mathbb{E}_0[\mathcal{M}_k] \leq T\varepsilon\sqrt{\frac{2}{\gamma_6}\mathbb{E}_0[\mathcal{N}_k]}.$$

*Proof.* A simple application of the Pinsker's inequality shows that

$$\mathbb{E}_k[\mathcal{M}_k] - \mathbb{E}_0[\mathcal{M}_k] \leq T\sqrt{\frac{1}{2}\text{KL}(\mathbb{P}_0, \mathbb{P}_k)}. \tag{20}$$

By Lemma 5.9 and by the standard KL decomposition theorem of the KL divergence [Cesa-Bianchi and Lugosi, 2006, Chapter 6], we have that the KL divergence between $\mathbb{P}_0$ and $\mathbb{P}_k$ only depends on the expected number of times that exploring actions were played, and on the KL divergence of the feedback distributions in such regions:

$$\text{KL}(\mathbb{P}_0, \mathbb{P}_k) = \mathbb{E}_0[\mathcal{N}_k] \cdot \text{KL}(\mathcal{H}_0, \mathcal{H}_k),$$

where, for each $k$, $\mathcal{H}_k$ is a discrete distribution on the 4 possible outcomes of the two-bit feedback. Formally,

$$\mathcal{H}_k(z) = \begin{cases} (k+1)\gamma_1 + \varepsilon\mathbb{I}\{k \neq 0\} + \gamma_6, & \text{if } z = (1,1) \\ \gamma_1 - \varepsilon\mathbb{I}\{k \neq 0\} + \gamma_1(N-2-k) + \gamma_6, & \text{if } z = (0,1) \\ (k+1)\gamma_1 - \varepsilon\mathbb{I}\{k \neq 0\} + \gamma_6 + \gamma_5 + \gamma_4(k+1) + \gamma_3^{\text{tot}}, & \text{if } z = (1,0) \\ \gamma_1 + \varepsilon\mathbb{I}\{k \neq 0\} + \gamma_1(N-2-k) + \gamma_6 + \gamma_4(N-k-1) & \text{if } z = (0,0) \end{cases}$$

By upperbounding the KL divergence with the $\chi^2$-distance [Lattimore and Szepesvári, 2020, Chapter 14], we immediately obtain that

$$\text{KL}(\mathcal{H}_0, \mathcal{H}_k) \leq \chi^2(\mathcal{H}_0, \mathcal{H}_k) = \sum_{z \in \{0,1\}^2} \frac{(\mathcal{H}_0(z) - \mathcal{H}_k(z))^2}{\mathcal{H}_0(z)} \leq \varepsilon^2\frac{4}{\gamma_6}, \tag{21}$$

where the last inequality holds since $\mathcal{H}_0(z) \geq \gamma_6$. $\qquad\square$

**Lemma 5.12.** *There is an instance $k \in \{0, \ldots, N-1\}$ and an absolute constant $c \in (0,1)$ such that:*

$$R_T^k \geq c \cdot \min\left(\frac{N}{\varepsilon^2}, \varepsilon T\right).$$

*Proof.* By summing over instances $k \in \{1, \ldots, N-1\}$ the result of Lemma 5.11 and using Jensen's inequality, we obtain:

$$\frac{1}{N-1}\sum_{k=1}^{N-1}(\mathbb{E}_k[\mathcal{M}_k] - \mathbb{E}_0[\mathcal{M}_k]) \leq \varepsilon T\sqrt{\frac{2}{\gamma_6}\frac{1}{N-1}\sum_{k=1}^{N-1}\mathbb{E}_0[\mathcal{N}_k]}. \tag{22}$$

Then, by rearranging Lemma 5.10 by and substituting $c_2 = c_1 - \gamma_5\frac{1-\ell}{2}$ we obtain:

$$\sum_{t=1}^{T}\mathbb{E}_k[\text{GFT}(p_t, q_t, s, b)] \leq \mathbb{E}_k\left[c_1 T + \rho \cdot \varepsilon\mathcal{M}_k - \gamma_5\frac{1-\ell}{2}\sum_{j=1}^{N-1}\mathcal{N}_j\right],$$

and by summing over $k \in \{1, \ldots, N-1\}$, dividing by $N-1$, and using Equation (22) we get

$$\frac{1}{N-1} \sum_{t=1}^{T} \sum_{k=1}^{N-1} \mathbb{E}_k[\mathrm{GFT}(p_t, q_t, s, b)] \le c_1 T + \rho\varepsilon \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_k[\mathcal{M}_k]$$

$$\le c_1 T + \rho\varepsilon \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{M}_k] + \rho\varepsilon^2 T \sqrt{\frac{2}{\gamma_6} \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{N}_k]}.$$

Then, the average of the regret $R_T^k$ over instances $k \in \{1, \ldots, N-1\}$ can be lower bounded by

$$\frac{1}{N-1} \sum_{k=1}^{N-1} R_T^k \ge T(c_1 + \rho\varepsilon) - \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_k[\mathrm{GFT}(p_t, q_t, s, b)]$$

$$\ge T(c_1 + \rho\varepsilon) - \left( c_1 T + \rho\varepsilon \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{M}_k] + \rho\varepsilon^2 T \sqrt{\frac{2}{\gamma_6} \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{N}_k]} \right)$$

$$\ge \rho\varepsilon T \left( \frac{1}{2} - \varepsilon \sqrt{\frac{2}{\gamma_6} \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{N}_k]} \right), \tag{23}$$

where the first inequality follows by Lemma 5.6 while the last inequality holds for any $N \ge 2$. Then, we divide the analysis in two cases. Intuitively, the first one correspond to the cases in which the algorithm does not explore enough (*i.e.*, $\sum_k \mathbb{E}_0[\mathcal{N}_k]$ is small) and, therefore, it cannot correctly identify the instance. In the second case the algorithm spends a large time exploring (*i.e.*, $\sum_k \mathbb{E}_0[\mathcal{N}_k]$ is large), and thereby accumulates large regret (by Lemma 5.8).

Formally, if $\varepsilon \sqrt{\frac{2}{\gamma_6} \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{N}_k]} \le \frac{1}{4}$ then Equation (23) implies that the average regret over instances $k \in \{1, \ldots, N-1\}$ is at least

$$\frac{1}{N-1} \sum_{k=1}^{N-1} R_T^k \ge \frac{1}{4}\rho\varepsilon T \ge \frac{1}{10^3}\varepsilon T.$$

Then, there must exist at least an instance $k \in \{1, \ldots, N-1\}$ in which $R_T^k = \Omega(\varepsilon T)$.

Otherwise, if $\varepsilon \sqrt{\frac{2}{\gamma_6} \frac{1}{N-1} \sum_{k=1}^{N-1} \mathbb{E}_0[\mathcal{N}_k]} \ge \frac{1}{4}$, then the regret of the base instance can be upper bounded by

$$R_T^0 \ge \frac{\gamma_5}{3} \mathbb{E}_0 \left[ \sum_{k=1}^{N-1} \mathcal{N}_k \right] \ge \frac{\gamma_5}{3} \frac{\gamma_6}{2} \left( \frac{1}{4\varepsilon} \right)^2 (N-1) \ge \frac{\gamma_5}{3} \frac{\gamma_6}{4} \left( \frac{1}{4\varepsilon} \right)^2 N \ge \frac{1}{10^6} \frac{N}{\varepsilon^2},$$

and hence in the base instance the regret is at least $\Omega\left( \frac{N}{\varepsilon^2} \right)$. □