

No-Regret Online Autobidding Algorithms in First-price Auctions

Yuan Deng^{*1}, Yilin Li^{†2}, Wei Tang^{‡2}, and Hanrui Zhang^{§2}

¹Google Research

²Chinese University of Hong Kong

Abstract

Automated bidding to optimize online advertising with various constraints, e.g. ROI constraints and budget constraints, is widely adopted by advertisers. A key challenge lies in designing algorithms for non-truthful mechanisms with ROI constraints. While prior work has addressed truthful auctions or non-truthful auctions with weaker benchmarks, this paper provides a significant improvement: We develop online bidding algorithms for repeated first-price auctions with ROI constraints, benchmarking against the optimal randomized strategy in hindsight. In the full feedback setting, where the maximum competing bid is observed, our algorithm achieves a near-optimal $\tilde{O}(\sqrt{T})$ regret bound, and in the bandit feedback setting (where the bidder only observes whether the bidder wins each auction), our algorithm attains $\tilde{O}(T^{3/4})$ regret bound.¹

1 Introduction

Autobidding has become the predominant paradigm in online advertising. The idea is that advertisers compete for ad opportunities through repeated auctions, where each advertiser employs an automated bidding algorithm (i.e., an autobidder) to bid on behalf of the advertiser. Autobidders are essentially online constrained optimizers: As auctions happen over time, an autobidder maximizes on the fly a certain objective quantity, e.g., clicks or conversions, subject to constraints, e.g., the total spending cannot exceed a predetermined budget, or the ratio between the return and the investment (i.e., the return-on-investment, or ROI ratio) must be at least a given target ratio. Normally, advertisers are free to customize these components of autobidders to better serve their own goals.

Over years, advertisers seem to have converged to the following common practice: An overwhelming majority of advertisers use autobidders that behave like *ROI-constrained value maximizers*. This highlights the importance of designing good online bidding algorithms in this very setting. The problem, roughly speaking, is the following: Initially, an auction mechanism (e.g., first-price or second-price) is chosen externally by the advertising platform, which is used in all subsequent auctions. At each step, an auction happens, and the algorithm observes the “value” of winning in

^{*}dengyuan@google.com

[†]ylli25@cse.cuhk.edu.hk

[‡]weitang@cuhk.edu.hk

[§]hanrui@cse.cuhk.edu.hk

¹A conference version of this work has been accepted in the proceeding of the 39th Conference on Neural Information Processing Systems (NeurIPS’25).

this auction. Then, the algorithm must immediately submit a bid based on historical observations and the value, after which the algorithm observes the outcome of the auction, i.e., whether it has won and possibly also the highest bid by other advertisers. In the long run, the algorithm must (approximately) maximize the total value obtained over all auctions, subject to the constraint that the overall ROI ratio is at least a predetermined target ratio. The key challenge here is that at each step, the algorithm must decide the bid without knowing the future, or even the competing bid in the current auction. Nonetheless, it must ensure that its decisions are almost optimal from hindsight (i.e., it has *no regret*), and the ROI constraint is approximately satisfied over the entire horizon.

Indeed, considerable effort has been invested into this task. Here, we refrain from a prolonged discussion (see Section 1.1 for further related work) and focus on the most related results. Feng et al. [19] study the problem when the auction mechanism is *truthful*, and give an almost optimal no-regret online bidding algorithm. Despite its undoubted importance, one severe limitation of their result is that it does not apply to some of the most common auction mechanisms in reality, including first-price auctions and generalized second-price auctions — for these mechanisms are not truthful. To address this issue, Aggarwal et al. [3] design alternative no-regret bidding algorithms that work for nontruthful auctions. However, their algorithms only have no regret against a weaker benchmark, i.e., the optimal *Lipschitz-continuous* bidding strategy from hindsight, which in general can be much worse than the unconditionally optimal bidding strategy. This leaves the following question open: *Is it possible to design a no-regret online bidding algorithm against the unconditionally optimal strategy for nontruthful auction mechanisms?*

Our results. We provide an affirmative answer to the above question. For concreteness, we derive our results for first-price auctions (which is the most representative nontruthful auction mechanism), but our results easily generalize to other auction mechanisms. We consider two natural feedback models commonly studied in the literature:

- *Full feedback*: the algorithm learns the auction outcome and the competing bid, after each auction;
- *Bandit (or one-bit) feedback*: the algorithm learns only the outcome of the auction and nothing else after each auction.

The latter setting is more challenging than the former, since the algorithm obtains less information.

In the full feedback setting, we present an online bidding algorithm that achieves the *optimal* (up to polylog factors) regret bound of $\tilde{O}(\sqrt{T})$ against the *unconditional, possibly randomized, optimal strategy* in hindsight, where T is the time horizon. This generalizes the result by Feng et al. [19] and strengthens the one by Aggarwal et al. [3] in the same setting. In the bandit feedback setting, we present an algorithm that achieves a regret bound of $\tilde{O}(T^{3/4})$ against the unconditionally optimal strategy, which strengthens the result by Aggarwal et al. [3] in the same setting. These results paint a considerably clearer picture of online bidding algorithms for ROI-constrained value maximizers.

Technique overview. The key technical challenge in designing online autobidding algorithms, which has also been observed by Aggarwal et al. [3], is that unconditionally optimal bidding strategies appear to have rich structures. This tends to baffle existing techniques, which normally require the class of strategies optimized over to be “simple” and / or “structured” in some way. Below we first briefly review how prior work deals with this challenge.

Feng et al. [19] restrict their attention to truthful auctions, in which optimal bidding strategies admit a clean characterization: Without loss of generality, they only need to consider “uniform bid scaling”, i.e., bidding strategies defined by a single multiplier θ , such that the strategy places a bid of

$\theta \cdot v$ when the value of winning is v . Restricted to such bidding strategies, one essentially only needs to optimize the multiplier θ , which allows Feng et al. [19] to essentially reduce their problem (in a nontrivial way) to one that can be handled within the powerful mirror descent framework. Aggarwal et al. [3] take an approach that is in a sense more general: They consider bidding strategies that are Lipschitz-continuous, which subsume uniform bid scaling as a subclass. By doing so, Aggarwal et al. [3] manage to further incorporate technical insights from the Lipschitz bandits literature, which help establish nontrivial regret bounds for nontruthful auctions.

We take an approach that is in spirit similar to [19]. We first make a more general structural observation, that even in nontruthful auctions, optimal (and possibly randomized) bidding strategies are *essentially parametrized by a single number*, which roughly corresponds to the marginal ROI of the bidding strategy. We present an algorithmically compatible characterization of this class of strategies (see Section 3), which allows us to (1) generalize the mirror-descent-based algorithm used in [19] to nontruthful auctions and (2) without loss of generality, focus on deterministic bidding strategies, provided one condition: The algorithm needs to *know the distribution of competing bids*.² This condition is crucial, since the structure of optimal bidding strategies necessarily depends on the distribution of competing bids, and without knowing the former, one simply cannot guarantee the conditions for the mirror-descent primitive to produce the desired regret bound.³ In our setting, of course the algorithm cannot know this distribution beforehand. Below we discuss how this requirement can be removed.

We first show that the requirement can be *relaxed* without hurting the regret too much. In particular, we show that if we run the generalized mirror-descent-based algorithm with a slightly inaccurate distribution of competing bids given as input, then its per-round regret can only degrade to the extent that the input distribution is inaccurate. With that as a primitive, we construct an algorithm that “bootstraps” itself without any prior knowledge about the environment in the full feedback setting (see Section 4). The algorithm “restarts” itself periodically as it gathers more information about the distribution of competing bids. Upon each restart, it runs a fresh instance of the mirror-descent-based primitive with the latest estimated distribution as its input, which allows it to gradually lower the per-round regret over time as the estimation becomes more accurate. By properly pacing the restarts, we manage to beat the $O(T^{2/3})$ bound via simple exploration-exploitation and achieve an almost optimal bound of $\tilde{O}(\sqrt{T})$.

In the bandit feedback setting, the above strategy does not work any more, and we have to fall back to simple exploration-exploitation. In our setting, even this conceptually simple high-level strategy requires nontrivial technical ingredients to implement: In fact, here we need a slightly stronger property of the primitive, because we can no longer guarantee a *uniform* estimation error of the distribution of competing bids. We show that the property holds for a slightly adapted version of the primitive used in the full feedback setting. With that as a building block, by optimizing the granularity of discretization and the tradeoff between exploration and exploitation, we manage to establish a regret bound of $\tilde{O}(T^{3/4})$ (see Section 5), which matches the guarantee established by Aggarwal et al. [3] against the weaker Lipschitz-continuous benchmark. Notably, the ideal $O(\sqrt{T})$ bound is impossible in this setting due to a lower bound of $\Omega(T^{2/3})$ by Aggarwal et al. [3]. We leave closing this regret gap as an important future direction.

²Note that the problem remains nontrivial, because the algorithm does not know the distribution of its value.

³This is not an issue in [19] since the structure they need is independent of the distribution of competing bids.

1.1 Related Work

Most closely related to our work is the line of research on online bidding algorithms in repeated auctions, with ROI-constrained value-maximizing bidders. In the paradigm with ROI constrained value maximizing bidders, Feng et al. [19] propose an algorithm based on the primal-dual framework proposed by Balseiro et al. [6], which achieves a nearly optimal $\tilde{O}(\sqrt{T})$ regret bound for truthful auctions. Castiglioni et al. [13] propose an algorithm that achieves no regret in repeated, possibly nontruthful, auctions, under the assumption that the spaces of values and bids are both finite. Aggarwal et al. [3] design an algorithm which works with nontruthful auctions and continuous value / bid space, and achieves no regret against any Lipschitz-continuous bidding strategy. Our work generalize / strengthen all these results by achieving the nearly optimal $\tilde{O}(\sqrt{T})$ regret against the unconditionally optimal bidding strategy without further assumptions.⁴ Also conceptually related is the work by Vijayan et al. [33], who consider a “harder” setting where the algorithm does not even observe the value of winning, and instead, must learn to estimate this value on the fly. Their result is orthogonal to ours, in particular because they also assume a finite space of bids.

Beyond autobidding, there is an extremely rich body of research on online bidding algorithms in repeated auctions. To name a few results: Balseiro et al. [6] propose a primal-dual framework and apply their framework to repeated second-price auctions with budget constraints and derive no-regret bidding algorithms. Wang et al. [34] study the same problem in repeated first-price auctions. Cesa-Bianchi et al. [14] study online bidding algorithms in first-price auctions without budget or ROI constraints, and pin down the optimal regret in a number of settings. Kumar et al. [26] present no-regret bidding algorithms that are strategically robust. Susan et al. [32] design no-regret bidding algorithms for multi-platform settings, where the bidding algorithm repeatedly participates in multiple parallel auctions. More generally, there is a long line of research on no-regret algorithms and dynamics for allocation problems under constraints, e.g., [5, 9, 11, 12, 20, 21, 22, 23, 24, 25, 29, 31, 35]. We refrain from an extensive discussion since these results concern fairly different problems from ours, both conceptually and technically.

Our results are also conceptually related to the growing literature on the design and analysis of auctions and marketplaces with autobidders. Prior research in this direction concerns various topics, including the efficiency of traditional auction mechanisms [2, 15, 16, 18, 27], the design of (approximately) optimal mechanisms [4, 8, 10, 17, 28, 30], etc. See the recent survey by Aggarwal et al. [1] for a comprehensive exposition.

2 Preliminaries

We consider an online bidding setting where a single bidder participates in single-item first-price auctions⁵, repeated in T rounds. At each time step $t \in [T]$, the learner realizes the value v_t in the current round and then submits a bid b_t based on the current value and the historical information up to now. The bidder wins the item if her submitted bid is not less than the (highest) competing bid d_t ($d_t \in [0, 1]$) in this round, and we denote by $x_t(b_t, d_t) = \mathbf{1}\{b_t \geq d_t\}$ the ex post allocation outcome of the current round. If the bidder wins the item, she pays a price b_t which equals to

⁴The $\tilde{O}(\sqrt{T})$ bound applies to the full feedback setting. In the bandit feedback setting, we match the bound by Aggarwal et al. [3] against a weaker benchmark, thereby also strengthening their result.

⁵We consider first-price auctions for simplicity. Our results generalize to other non-truthful auctions: The main technical difference would be in our Theorem 3.1, where the construction of F_{conv} should be adapted in accordance with the auction mechanism.

her submitted bid, otherwise she pays 0. Therefore, the payment function, denoted by $p_t(b_t, d_t)$, is $p_t(b_t, d_t) = b_t \cdot \mathbf{1}\{b_t \geq d_t\}$. We focus on value-maximizing bidder where the bidder's utility at time t is given by $r_t(b_t, d_t) = v_t \cdot x_t(b_t, d_t)$.

Benchmark and regret definitions. We assume that the sequence of value v_t is drawn independently and identically (i.i.d.) from an unknown distribution, whose CDF is denoted as H . The bidder's objective is to design an online bidding algorithm ALG that maximizes her total realized utility subject to a Return-On-Investment (ROI) constraint.⁶ Namely, her total utility must be at least ρ times her total payment: $\rho \cdot \sum_{t \in [T]} p_t(b_t, d_t) \leq \sum_{t \in [T]} v_t \cdot x_t(b_t, d_t)$ where $\rho > 0$ is the target ROI ratio of the bidder.

$$\begin{aligned} \max \quad & \mathbb{E} \left[\sum_{t \in [T]} v_t \cdot x_t(b_t, d_t) \right] \\ \text{s.t.} \quad & \mathbb{E} \left[\rho \cdot \sum_{t \in [T]} p_t(b_t, d_t) \right] \leq \mathbb{E} \left[\sum_{t \in [T]} v_t \cdot x_t(b_t, d_t) \right], \end{aligned} \tag{\mathcal{P}}$$

where the expectation is over the possible randomness of the bid sequence $\{b_t\}$ and the randomness of the maximum competing bid $\{d_t\}$. W.l.o.g., we assume that $\rho = 1$.⁷

We assume that the (maximum) competing bid d_t at each round is i.i.d. realized from an unknown competing bid distribution F . In this context, the allocation and the bidder's payment in round t are determined by her submitted bid b_t : the allocation function $x_t(b_t, d_t) = F(b_t)$ and the payment function $p_t(b_t, d_t) = b_t \cdot F(b_t)$.

We denote the bidder's value sequence by $V^T = (v_t)_{t \in [T]}$. An online bidding algorithm takes input of the bidder's historical information and the realized value at the current round, and outputs a (possibly randomized) bid. Given an online algorithm ALG , let $\text{Reward}(\text{ALG} | (F, V^T))$ (resp. $\text{Payment}(\text{ALG} | (F, V^T))$) denote the bidder's total expected utility (resp. total expected payment) given the value sequence V^T and competing bid distribution F :

$$\begin{aligned} \text{Reward}(\text{ALG} | (F, V^T)) &= \mathbb{E}_{(b_t) \sim \text{ALG}} \left[\sum_{t \in [T]} v_t \cdot F(b_t) \right]; \\ \text{Payment}(\text{ALG} | (F, V^T)) &= \mathbb{E}_{(b_t) \sim \text{ALG}} \left[\sum_{t \in [T]} b_t \cdot F(b_t) \right]. \end{aligned}$$

Thus, the ROI violation is given by $\Delta(\text{ALG} | (F, V^T)) = \text{Payment}(\text{ALG} | (F, V^T)) - \text{Reward}(\text{ALG} | (F, V^T))$.

Let OPT be the corresponding optimal bidding algorithm in hindsight for the program \mathcal{P} , i.e., OPT has knowledge about the competing bid distribution F and the realized value sequence V^T . Thus, an algorithm's regret relative to OPT under the input value sequence V^T is

$$\text{Regret}(\text{ALG} | (F, V^T)) = \text{Reward}(\text{OPT} | (F, V^T)) - \text{Reward}(\text{ALG} | (F, V^T)).$$

We measure the performance of an algorithm ALG by its expected regret over the distributions F and H , defined as follows:

$$\text{Regret}(\text{ALG} | (F, H)) = \mathbb{E}_{V^T \sim H^T} [\text{Regret}(\text{ALG} | (F, V^T))]. \tag{1}$$

⁶For pedagogical reasons, we focus exclusively on ROI constraints. Our results and analysis can be readily generalized to the setting with additional budget constraints, e.g., through techniques introduced in [7].

⁷For $\rho \neq 1$, we can rescale bidder's value to be $v_t \cdot \rho$, and all our analysis/results can be carried over similarly.

and expected ROI violation over the distribution F and H :

$$\Delta(\text{ALG} | (F, H)) = \mathbb{E}_{V^T \sim H^T} [\Delta(\text{ALG} | (F, V^T))] .$$

The feedback models. We now specify the feedback that the bidder receives at the end of each round. In this paper, we consider the following two feedback models.

- **Full feedback:** At the end of round t , the bidder observes the competing bid d_t .
- **Bandit feedback:** At the end of round t , the bidder only receives a binary signal $\text{win}_t \in \{0, 1\}$ indicating whether she won the item or not, i.e., $\text{win}_t = \mathbf{1}\{d_t \leq b_t\}$.

3 Optimal Randomized Autobidding Strategy

In this section, we first present the bidder’s optimal bidding strategy in hindsight, where the sequence of values V^T is known. We then design an online bidding algorithm that has low regret and low ROI violation when competing bid distribution F is known, while value distribution H remains unknown.

We first define the following allocation-payment curve, which will be helpful for our analysis.

Definition 3.1 (Allocation-payment curve G). *Given a competing bid distribution $F \in \Delta([0, 1])$, an allocation-payment curve $G : [0, 1] \rightarrow [0, 1]$ is defined as: $G(b \cdot F(b)) = F(b)$ for all $b \in [0, 1]$.*

The allocation-payment curve G captures the relationship between the payment and the allocation (as well as the reward, which is proportional to the allocation when the value of winning is fixed) induced by any bid b .

3.1 Characterizing Optimal Bidding Strategy

We start with the following example which shows the suboptimality of deterministic bidding strategies.

Example 3.2 (Suboptimality of deterministic strategy). *Consider following competing bid distribution: $\Pr[d=0] = \Pr[d=1] = 1/2$. The bidder always has a value $v = 1/2$. Let $T = 1$, and the target ROI ratio $\rho = 1$. One optimal deterministic bidding strategy would be to always submit a bid $b = 0$, which yields an expected reward $1/4$ and an expected payment 0 — note that restricted to deterministic bidding, there is no way to utilize this “slackness” in the ROI constraint.⁸ On the other hand, consider the following randomized bidding strategy: $\Pr[b=0] = 2/3, \Pr[b=1] = 1/3$. This strategy yields an expected reward $1/3$ and an expected payment $1/3$. Both strategies satisfy the ROI constraint, while the randomized strategy gives a strictly higher reward.*

Definition 3.2 shows that to maximize the bidder’s expected reward, it may be necessary to solve program \mathcal{P} over all randomized bidding strategies. This is (at least superficially) incompatible with the primal-dual framework that underlies all previous results on online autobidding, which presents new technical challenges. To address this issue, we present a constructive reduction that brings us back to the deterministic world, where we can cast the primal-dual framework more naturally. In

⁸Here and in the rest of the paper, we assume ties are broken in favor of the bidding algorithm. This makes no difference when the distribution of competing bids is continuous.

particular, we show that given any input value sequence V^T and any competing bid distribution F , it is always possible to construct another distribution F_{conv} that captures essentially the same tradeoff between reward and payment as F induces, with one key difference: With F replaced by F_{conv} , we can focus our attention to deterministic bidding strategies without loss of generality.

Theorem 3.1 (Optimal randomized bidding strategy). *Fixing any input value sequence V^T . For any competing bid distribution F , there exists a distribution $F_{\text{conv}} \in \Delta([0, 1])$ such that: For any (randomized) strategy under V^T and F , there is a deterministic strategy under V^T and F_{conv} inducing reward and payment that are both no worse, and vice versa.*

In the proof of Theorem 3.1, we explicitly construct such distribution F_{conv} : the distribution F_{conv} is constructed such that its allocation-payment curve, denoted by G_{conv} , is *essentially* the concave envelope of the allocation-payment curve G induced by the original competing bid distribution F . With Theorem 3.1, we can simplify our algorithm design as follows: instead of designing an online bidding algorithm that has low regret against the optimal randomized bidding strategy under (F, V^T) , we can turn to designing an algorithm that has low regret against the optimal deterministic strategy under the corresponding (F_{conv}, V^T) .

We conclude this subsection with a proof sketch of Theorem 3.1. All the missing proofs for the technique lemmas are deferred to Section A.

Proof of Theorem 3.1. Fixing any input value sequence V^T . For any competing bid distribution F , the proof of Theorem 3.1 consists of three main steps:

- **Step 1 – Prove that the “concave envelope” G_{conv} ⁹ of allocation-payment curve G attains the maximum expected reward (see Theorem A.2):** In our first step, we show that given any fixed value v , any (possibly randomized) bidding strategy that leads to an expected payment of x has an expected reward at most $v \cdot G_{\text{conv}}(x)$.
- **Step 2 – Construct the distribution F_{conv} (see Theorem A.3):** In this step, we construct such distribution F_{conv} that corresponds to an allocation-payment curve G_{conv} , which is essentially the concave envelope of the allocation-payment curve G induced by F .
- **Step 3 – Prove the reward equivalence (see Theorem A.4):** In this step, we show that any reasonable randomized bidding strategy under (F, V^T) is equivalent to generating a bid on the allocation-payment curve G_{conv} , which corresponds to a deterministic bidding strategy under (F_{conv}, V^T) , and vice versa. \square

3.2 No-Regret Autobidding Algorithm for Known Competing Bid Distribution

In this subsection, we describe a no-regret online bidding algorithm (see Algorithm 1) when the competing bid distribution F is known to the bidder, while the value distribution H still remains unknown. This algorithm will serve as a building block for our algorithm design when the both F and H are unknown to the bidder.

We first note that after introducing the Lagrangian multiplier $\lambda \geq 0$, we can reformulate the program \mathcal{P} as follows

$$\max \mathbb{E} \left[\sum_{t \in [T]} v_t \cdot F(b_t) + \min_{\lambda \geq 0} \lambda \cdot \sum_{t \in [T]} (v_t \cdot F(b_t) - b_t \cdot F(b_t)) \right]. \quad (2)$$

⁹Technically, G_{conv} might differ from the actual concave envelope of G , but this can happen only on the part of G_{conv} that never matters for our purposes. See the proof of Theorem A.3 for details.

Following a similar dual-prime framework proposed by [6], we proceed by choosing a bid b_t at round t to maximize the objective in Eqn. (2) using a fixed Lagrangian multiplier λ_t and then update λ_t to λ_{t+1} use the observed feedback. Notice that given a fixed λ_t at round t , the problem in Eqn. (2) is equivalent to

$$\max \mathbb{E} \left[\frac{1 + \lambda_t}{\lambda_t} \cdot v_t \cdot F(b) - b \cdot F(b) \right]. \quad (3)$$

As we show in Definition 3.2, the optimal bidding strategy may be randomized. Thus, using Theorem 3.1, instead of directly solving problem Eqn. (3), we consider solving a bid \tilde{b}_t that maximizes $\frac{1+\lambda_t}{\lambda_t} \cdot v_t \cdot F_{\text{conv}}(b) - b \cdot F_{\text{conv}}(b)$. We then use Theorem A.4 to construct the equivalent randomized bidding strategy to be implemented in practice. We then update the dual variable λ_t using the a generalized negative entropy in the mirror descent framework to adjust the penalty for the future rounds. We summarize the algorithm details in Algorithm 1.

Algorithm 1 No-regret bidding algorithm for known F

Input: Time horizon T , competing bid distribution F , the value sequence V^T (revealed sequentially).

Initialize: Initial dual variable $\lambda_1 = 1$ and the dual mirror descent step size $\alpha = 1/\sqrt{T}$

- 1: Let F_{conv} be the distribution defined as in Theorem 3.1 (its construction is in Theorem A.3).
 - 2: **for** $t \leftarrow 1$ to T **do**
 - 3: Observe the value v_t .
 - 4: Compute $\tilde{b}_t \leftarrow \operatorname{argmax}_{b \in [0,1]} \left[\frac{1+\lambda_t}{\lambda_t} \cdot v_t \cdot F_{\text{conv}}(b) - b \cdot F_{\text{conv}}(b) \right]$
 - 5: Given \tilde{b}_t , construct the equivalent randomized bidding strategy mentioned in Theorem 3.1 and realizes a bid b_t . /* Such randomized bidding strategy can be constructed efficiently, see Theorem A.4. */
 - 6: Submit the bid b_t (to the outer environment or algorithm).
 - 7: Update $g_t(\tilde{b}_t) = v_t \cdot F_{\text{conv}}(\tilde{b}_t) - \tilde{b}_t \cdot F_{\text{conv}}(\tilde{b}_t)$.
 - 8: Update $\lambda_{t+1} = \lambda_t \cdot \exp[-\alpha \cdot g_t(\tilde{b}_t)]$.
 - 9: **end for**
 - 10: **return** The bid sequence $\{b_t\}_{t \in [T]}$.
-

The performance of Algorithm 1 is detailed below, whose analysis follows similarly to [6, 19]. We defer the proof to Section A.

Proposition 3.2. *For any unknown value distribution H , the expected regret of Algorithm 1 under any competing bid distribution F is bounded by $\text{Regret}(\text{Algorithm 1} | (F, H)) = O(\sqrt{T})$, and the expected ROI violation is bounded by $\Delta(\text{Algorithm 1} | (F, H)) \leq 2 \log T \sqrt{T}$.*

4 No-Regret Autobidding Algorithm with Full Feedback

In Section 3, we present an algorithm (see Algorithm 1) when the competing bid distribution is known that achieves low expected regret and ROI violation. In practice, both the competing bid distribution F and the value distribution H may be unknown, and the bidder needs to use the observed feedback to learn these uncertainties by exploring different bids. In this section, we consider a full feedback model where the bidder can observe the realized competing bid d_t at the end of each round, and we provide an algorithm (see Algorithm 2) with low expected regret and low ROI violation.

Algorithm descriptions. At a high-level, Algorithm 2 operates as a stage-based algorithm. In each stage $m \in [M]$ where M is the number of stages determined shortly, we use the competing bids observed in all previous stages to maintain an empirical CDF \hat{F}_m for the underlying competing bid distribution F . Let $T_m \in \mathbb{N}^+$ be the number of rounds in stage m and $N_m = \sum_{m' \in [m-1]} T_{m'}$ be the number of total rounds right before the stage m . We maintain \hat{F}_m as follows:

$$\hat{F}_m(b) = \frac{1}{N_m} \sum_{i \in [N_m]} \mathbf{1}\{b \geq d_i\}, \quad b \in [0, 1], \quad (4)$$

To account for the estimation error and maintain optimism in our bidding strategy, we further define the optimistic CDF \bar{F}_m as

$$\bar{F}_m(b) = (\hat{F}_m(b) + \varepsilon_m) \wedge 1, \quad b \in [0, 1]. \quad (5)$$

where ε_m is a parameter carefully chosen to balance the exploration and exploitation tradeoff.¹⁰ We then implement Algorithm 1 by feeding the distribution \bar{F} (while the actual competing bids are still generated according to the true unknown competing bid distribution F). We summarize the algorithm details in Algorithm 2.

Algorithm 2 No-regret bidding algorithm with full feedback

Input: Time horizon T , the value sequence V^T (revealed sequentially).

Input: Let $M = \lceil \log(T+1) \rceil$ be the number of stages.

Initialize: Let \bar{F}_1 be an arbitrary distribution. $N_1 \leftarrow 0$

- 1: **for** $m \leftarrow 1$ to M **do**
 - 2: Current stage length: $T_m \leftarrow 2^{m-1}$.
 - 3: Feed the Algorithm 1 with inputs: $T \leftarrow T_m, F \leftarrow \bar{F}_m$. /* Note that we compute the submitted bid b_t and update dual variable λ_t from Algorithm 1 based on \bar{F}_m , while the actual competing bids are generated according to F . */
 - 4: $N_{m+1} \leftarrow T_m + N_m$
 - 5: Optimism parameter $\varepsilon_{m+1} \leftarrow \frac{\log(N_{m+1})}{\sqrt{N_{m+1}}}$.
 - 6: Update the empirical CDF \hat{F}_m to be \hat{F}_{m+1} according to Eqn. (4).
 - 7: Obtain \bar{F}_{m+1} defined as in Eqn. (5)
 - 8: **end for**
 - 9: **return** The bid sequence $\{b_t\}_{t \in [T]}$.
-

Algorithm performance. We summarize the performance of Algorithm 2 as follows.

Theorem 4.1. *For any unknown competing bid distribution F and unknown value distribution H , under the setting with full feedback, the expected regret of Algorithm 2 is bounded by $\text{Regret}(\text{Algorithm 2} | (F, H)) = \tilde{O}(\sqrt{T})$, and the expected ROI violation is bounded by $\Delta(\text{Algorithm 2} | (F, H)) = \tilde{O}(\sqrt{T})$.*

We conclude this section by providing a proof sketch for Theorem 4.1. All missing proofs of the technical lemmas are deferred to Section B.

¹⁰As a shorthand, we write $x \wedge y := \min\{x, y\}$.

Proof of Theorem 4.1. We first discuss the regret and ROI violation incurred with any fixed value sequence V^T . The analysis of regret and ROI violation of algorithm 2 consists of three main steps:

- **Step 1 – Accuracy of the optimistic CDF \bar{F} (see Theorem B.1):** In this step, we proved that when $\varepsilon_m = \Theta\left(\frac{\log N_m}{\sqrt{N_m}}\right)$, the CDF \bar{F}_m (defined in Eqn. (5)) serves as an optimistic estimate to the true underlying competing bid distribution ($\bar{F}_m \geq F$ with high probability), while still remains close, specifically we have $\sup_b |\bar{F}_m(b) - F(b)| \leq 2\varepsilon_m$ with high probability.
- **Step 2 – The performance approximation when implementing Algorithm 1 with estimated distributions (see Theorem B.2):** In this step, we establish the performance guarantee (evaluating in an environment with competing bid distribution F) of the bids generated by feeding Algorithm 1 with a distribution F^\dagger that satisfies $F^\dagger \geq F$ and $\sup_b |F_m^\dagger(b) - F(b)| \leq 2\varepsilon$,

$$\text{Regret}(\text{ALG1}(F^\dagger) \mid (F, V^T)) = O(\sqrt{T}) + 2\varepsilon T, \quad \Delta(\text{ALG1}(F^\dagger) \mid (F, V^T)) = \tilde{O}(\sqrt{T}) + 2\varepsilon T.$$

Here we denote by $\text{ALG1}(F^\dagger)$ the bids generated by feeding Algorithm 1 with a distribution F^\dagger . In stage m , when we use the estimated optimistic \bar{F}_m as input for Algorithm 1, Theorem B.2 demonstrates that both the regret and ROI violation will be amplified by $T \cdot 2\varepsilon_m$ compared to using the true CDF F as input.

- **Step 3 – Aggregating the regret and ROI violation over M stages (see Theorem B.5 and Theorem B.6):** With ε_m selected in Step 1, Theorem B.5 proved that the regret of Algorithm 2 is bounded by $\tilde{O}(\sqrt{T})$. Theorem B.6 proved that the ROI violation is bounded by $\tilde{O}(\sqrt{T})$.

Thus, when a value sequence V^T is i.i.d. drawn from H^T , the expected regret and ROI violation over the distribution F and H satisfies

$$\text{Regret}(\text{Algorithm 2} \mid (F, H)) = \mathbb{E}_{V^T \sim H^T} [\text{Regret}(\text{Algorithm 2} \mid (F, V^T))] = \tilde{O}(\sqrt{T}).$$

$$\Delta(\text{Algorithm 2} \mid (F, H)) = \mathbb{E}_{V^T \sim H^T} [\Delta(\text{Algorithm 2} \mid (F, V^T))] = \tilde{O}(\sqrt{T}). \quad \square$$

5 No-Regret Autobidding Algorithm with Bandit Feedback

In this section, we consider a more challenging learning setting in which the bidder receives only bandit feedback at the end of each round – specifically, whether they won the item or not. This bandit feedback model is common in practice and has been extensively studied in the literature (see, e.g., [3, 14]). The main results in this section are an online autobidding algorithm (see Algorithm 3) with the provable performance (see Theorem 5.1) under such bandit feedback.

Algorithm descriptions. At a high-level, our proposed Algorithm 3 operates as an explore-then-exploit type algorithm. The algorithm starts with an exploration stage that explores different bids to obtain sufficient knowledge about the underlying competing bid distribution. Then, the algorithm implements a “conservative” variant of Algorithm 1 with the estimated competing bid distribution over the remaining rounds.

In the exploration phase, to explore different bids, we discretize the bid range $[0, 1]$ into K equal intervals with grid points $\{c_k\}_{k \in [K] \cup \{0\}}$, where $c_k = k/K$ for $k \in [K] \cup \{0\}$. We then bid each grid point c_k consecutively for M rounds, and then use the observed bidding outcomes to compute an

empirical estimate \hat{F} for the underlying competing bid distribution F : for $\ell \in [1 : K - 1]$

$$\hat{F}(c_\ell) = \max_{k \in [\ell] \cup \{0\}} \frac{1}{M} \cdot \sum_{t \in [k \cdot M + 1 : (k+1) \cdot M]} \mathbf{1}\{b_t \geq d_t\} , \quad (6)$$

and we let $\hat{F}(c_K) = 1$. We also extend the definition of \hat{F} to any bid $b \in [0, 1]$ as a step function: $\hat{F}(b) = \hat{F}(c_k)$ if $b \in [c_k, c_{k+1})$. After exploring all the grid bids, we can have sufficient observations to maintain a good estimate \hat{F} for the underlying the competing bid distribution. Similar to Algorithm 2, we also define the optimistic CDF \bar{F} as follows:

$$\bar{F}(b) = (\hat{F}(c_k) + \varepsilon) \wedge 1, \quad \text{if } b \in [c_k, c_{k+1}) , \quad (7)$$

where ε again is a parameter carefully chosen to balance the exploration and exploitation tradeoff.

In Algorithm 2 with full feedback, after obtaining the optimistic CDF \bar{F} , we directly implement Algorithm 1 assuming the underlying competing bid distribution is \bar{F} . However, this approach may fail under our bandit feedback setting, as we can no longer guarantee a uniform error bound on \bar{F} (i.e., $\sup_b |F(b) - \bar{F}(b)|$ might be large). Instead, in the exploitation phase of Algorithm 3, we implement Algorithm 1 for the remaining rounds in a “conservative” way. In particular, we define a conservative version of the optimistic CDF \bar{F}^c as follows:

$$\bar{F}^c(b) = \bar{F}((b + 1/K) \wedge 1), \quad b \in [0, 1] . \quad (8)$$

We then feed Algorithm 1 with the input $T - K \cdot M$ as the time horizon, and \bar{F}^c as if the underlying competing bid distribution is \bar{F}^c to generate a sequence of bids $\{b_t^c\}$ for the remaining rounds. Instead of directly using these bids, we submit the bids also in a conservative way. In particular, we bid $\{b_t\}$ where each bid $b_t \leftarrow (b_t^c + 1/K) \wedge 1$.

Algorithm 3 No-regret bidding algorithm with bandit feedback

Input: Time horizon T , the value sequence V^T (revealed sequentially).

Initialize: Initialize dual variable $\lambda_1 = 1, K, M$.

- ```

/* Exploration phase */
1: Discretize $[0, 1]$ to obtain grid points $\{c_0, c_1, c_2, \dots, c_K\}$.
2: for $k \leftarrow 0$ to $K - 1$ do
3: Bid c_k consecutively for M time rounds.
4: Observe the allocation outcomes $\mathbf{1}\{b_t \geq d_t\}_{t \in [k \cdot M + 1 : (k+1) \cdot M]}$.
5: end for
6: Obtain the empirical CDF \hat{F} , the optimistic CDF \bar{F} and the conservative CDF \bar{F}^c defined as in
Eqn. (6), (7) and (8), respectively.
/* Exploitation phase */
7: Feed the Algorithm 1 with inputs: $T \leftarrow T - K \cdot M, F \leftarrow \bar{F}^c$.
8: Let $\{b_t^c\}$ be the bids generated from Algorithm 1 with these input specifications.
9: Submit the bids $\{b_t\}$ over remaining rounds where each bid $b_t \leftarrow (b_t^c + 1/K) \wedge 1$.
10: return The bid sequence $\{b_t\}_{t \in [T]}$

```
- 

**Algorithm performance.** The performance of Algorithm 3 is summarized as follows:

**Theorem 5.1.** *For any unknown competing bid distribution  $F$  and unknown value distribution  $H$ , under the bandit feedback, the expected regret of Algorithm 3 with  $M = \Theta(\sqrt{T})$ ,  $K = \Theta(T^{1/4})$  is bounded by  $\text{Regret}(\text{Algorithm 3} | (F, H)) = \tilde{O}(T^{3/4})$ , and the expected ROI violation is bounded by  $\Delta(\text{Algorithm 3} | (F, H)) = \tilde{O}(T^{3/4})$ .*

The proof of Theorem 5.1 and all relevant lemmas (with proofs) are deferred to Section C.

## References

- [1] Gagan Aggarwal, Ashwinkumar Badanidiyuru, Santiago R Balseiro, Kshipra Bhawalkar, Yuan Deng, Zhe Feng, Gagan Goel, Christopher Liaw, Haihao Lu, Mohammad Mahdian, et al. 2024. Auto-bidding and auctions in online advertising: A survey. *ACM SIGecom Exchanges* 22, 1 (2024), 159–183.
- [2] Gagan Aggarwal, Ashwinkumar Badanidiyuru, and Aranyak Mehta. 2019. Autobidding with constraints. In *Web and Internet Economics: 15th International Conference, WINE 2019, New York, NY, USA, December 10–12, 2019, Proceedings* 15. Springer, 17–30.
- [3] Gagan Aggarwal, Giannis Fikoris, and Mingfei Zhao. 2025. No-regret algorithms in non-truthful auctions with budget and roi constraints. In *Proceedings of the ACM on Web Conference 2025*. 1398–1415.
- [4] Santiago Balseiro, Yuan Deng, Jieming Mao, Vahab Mirrokni, and Song Zuo. 2024. Optimal Mechanisms for a Value Maximizer: The Futility of Screening Targets. In *Proceedings of the 25th ACM Conference on Economics and Computation*. 1101–1101.
- [5] Santiago Balseiro, Negin Golrezaei, Mohammad Mahdian, Vahab Mirrokni, and Jon Schneider. 2019. Contextual bandits with cross-learning. *Advances in Neural Information Processing Systems* 32 (2019).
- [6] Santiago Balseiro, Haihao Lu, and Vahab Mirrokni. 2020. Dual mirror descent for online allocation problems. In *International Conference on Machine Learning*. PMLR, 613–628.
- [7] Santiago R Balseiro, Kshipra Bhawalkar, Zhe Feng, Haihao Lu, Vahab Mirrokni, Balasubramanian Sivan, and Di Wang. 2024. A Field Guide for Pacing Budget and ROS Constraints. In *International Conference on Machine Learning*. PMLR, 2607–2638.
- [8] Santiago R Balseiro, Yuan Deng, Jieming Mao, Vahab S Mirrokni, and Song Zuo. 2021. The landscape of auto-bidding auctions: Value versus utility maximization. In *Proceedings of the 22nd ACM Conference on Economics and Computation*. 132–133.
- [9] Santiago R Balseiro and Yonatan Gur. 2019. Learning in repeated auctions with budgets: Regret minimization and equilibrium. *Management Science* 65, 9 (2019), 3952–3968.
- [10] Xiaohui Bei, Pinyan Lu, Zhiqi Wang, Tao Xiao, and Xiang Yan. 2025. Optimal Auction Design for Mixed Bidders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 39. 13614–13621.
- [11] Christian Borgs, Jennifer Chayes, Nicole Immorlica, Kamal Jain, Omid Etesami, and Mohammad Mahdian. 2007. Dynamics of bid optimization in online advertisement auctions. In *Proceedings of the 16th international conference on World Wide Web*. 531–540.
- [12] Matteo Castiglioni, Andrea Celli, and Christian Kroer. 2022. Online learning with knapsacks: the best of both worlds. In *International Conference on Machine Learning*. PMLR, 2767–2783.
- [13] Matteo Castiglioni, Andrea Celli, and Christian Kroer. 2024. Online Learning under Budget and ROI Constraints via Weak Adaptivity. In *International Conference on Machine Learning*. PMLR, 5792–5816.

- [14] Nicolò Cesa-Bianchi, Tommaso Cesari, Roberto Colombari, Federico Fusco, and Stefano Leonardi. 2024. The role of transparency in repeated first-price auctions with unknown valuations. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*. 225–236.
- [15] Yuan Deng, Mohammad Mahdian, Jieming Mao, Vahab Mirrokni, Hanrui Zhang, and Song Zuo. 2024. Efficiency of the generalized second-price auction for value maximizers. In *Proceedings of the ACM Web Conference 2024*. 46–56.
- [16] Yuan Deng, Jieming Mao, Vahab Mirrokni, Hanrui Zhang, and Song Zuo. 2024. Efficiency of the first-price auction in the autobidding world. *Advances in Neural Information Processing Systems* 37 (2024), 139270–139293.
- [17] Yuan Deng, Vahab Mirrokni, and Hanrui Zhang. 2022. Posted pricing and dynamic prior-independent mechanisms with value maximizers. *Advances in Neural Information Processing Systems* 35 (2022), 24158–24169.
- [18] Yiding Feng, Brendan Lucier, and Aleksandrs Slivkins. 2024. Strategic budget selection in a competitive autobidding world. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*. 213–224.
- [19] Zhe Feng, Swati Padmanabhan, and Di Wang. 2023. Online Bidding Algorithms for Return-on-Spend Constrained Advertisers. In *Proceedings of the ACM Web Conference 2023*. 3550–3560.
- [20] Giannis Fikioris and Éva Tardos. 2023. Approximately stationary bandits with knapsacks. In *The Thirty Sixth Annual Conference on Learning Theory*. PMLR, 3758–3782.
- [21] Giannis Fikioris and Éva Tardos. 2023. Liquid welfare guarantees for no-regret learning in sequential budgeted auctions. In *Proceedings of the 24th ACM Conference on Economics and Computation*. 678–698.
- [22] Jason Gaitonde, Yingkai Li, Bar Light, Brendan Lucier, and Aleksandrs Slivkins. 2023. Budget Pacing in Repeated Auctions: Regret and Efficiency Without Convergence. In *14th Innovations in Theoretical Computer Science Conference (ITCS 2023)*, Vol. 251. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 52.
- [23] Yanjun Han, Zhengyuan Zhou, Aaron Flores, Erik Ordentlich, and Tsachy Weissman. 2020. Learning to bid optimally and efficiently in adversarial first-price auctions. *arXiv preprint arXiv:2007.04568* (2020).
- [24] Nicole Immorlica, Karthik Sankararaman, Robert Schapire, and Aleksandrs Slivkins. 2022. Adversarial bandits with knapsacks. *J. ACM* 69, 6 (2022), 1–47.
- [25] Thomas Kesselheim and Sahil Singla. 2020. Online learning with vector costs and bandits with knapsacks. In *Conference on Learning Theory*. PMLR, 2286–2305.
- [26] Rachitesh Kumar, Jon Schneider, and Balasubramanian Sivan. 2024. Strategically-Robust Learning Algorithms for Bidding in First-Price Auctions. In *Proceedings of the 25th ACM Conference on Economics and Computation*. 893–893.

- [27] Christopher Liaw, Aranyak Mehta, and Andres Perlroth. 2023. Efficiency of non-truthful auctions in auto-bidding: The power of randomization. In *Proceedings of the ACM Web Conference 2023*. 3561–3571.
- [28] Pinyan Lu, Chenyang Xu, and Ruilong Zhang. 2023. Auction design for value maximizers with budget and return-on-spend constraints. In *International Conference on Web and Internet Economics*. Springer, 474–491.
- [29] Brendan Lucier, Sarath Pattathil, Aleksandrs Slivkins, and Mengxiao Zhang. 2024. Autobidders with budget and roi constraints: Efficiency, regret, and pacing dynamics. In *The Thirty Seventh Annual Conference on Learning Theory*. PMLR, 3642–3643.
- [30] Hongtao Lv, Xiaohui Bei, Zhenzhe Zheng, and Fan Wu. 2023. Auction design for bidders with ex post roi constraints. In *International Conference on Web and Internet Economics*. Springer, 492–508.
- [31] Thomas Nedelec, Clément Calauzènes, Noureddine El Karoui, Vianney Perchet, et al. 2022. Learning in repeated auctions. *Foundations and Trends® in Machine Learning* 15, 3 (2022), 176–334.
- [32] Fransisca Susan, Negin Golrezaei, and Okke Schrijvers. 2023. Multi-platform budget management in ad markets with non-ic auctions. *arXiv preprint arXiv:2306.07352* (2023).
- [33] Sushant Vijayan, Zhe Feng, Swati Padmanabhan, Karthikeyan Shanmugam, Arun Suggala, and Di Wang. 2025. Online Bidding under RoS Constraints without Knowing the Value. In *Proceedings of the ACM on Web Conference 2025*. 3096–3107.
- [34] Qian Wang, Zongjun Yang, Xiaotie Deng, and Yuqing Kong. 2023. Learning to bid in repeated first-price auctions with budgets. In *International Conference on Machine Learning*. PMLR, 36494–36513.
- [35] Jonathan Weed, Vianney Perchet, and Philippe Rigollet. 2016. Online learning in repeated auctions. In *Conference on Learning Theory*. PMLR, 1562–1583.

## A Missing Proofs in Section 3

### A.1 Analysis of Randomized Autobidding Strategy

**Lemma A.1.** *For any  $F$ , the allocation-payment curve  $G$  is non-decreasing.*

*Proof of Theorem A.1.* Since  $b \geq 0$ , we have that  $b \cdot F(b)$  remains non-decreasing in  $b$ . Concretely, consider  $G(x_1) = F(b_1) \leq F(b_2) = G(x_2)$ . Due to the monotonicity of  $F(b)$ , we have  $b_1 < b_2$ . Therefore,  $x_1 = b_1 \cdot F(b_1) < b_2 \cdot F(b_2) = x_2$ . Hence, the curve  $G$  is a non-decreasing function.  $\square$

**Lemma A.2.** *Given any competing bid distribution  $F$ , let  $G$  be its allocation-payment curve and let  $G_{\text{concave}}(x)$  be its concave envelope. Fix any given value  $v \in [0, 1]$ . For any (possibly randomized) bidding strategy that leads to an expected payment of  $x$ , then its expected reward is at most  $v \cdot G_{\text{concave}}(x)$ .*

*Proof of Theorem A.2.* Let us fix an arbitrary competing bid distribution  $F$ . Fix a value  $v$ , and consider any randomized combination of bids  $\{b_i\}_{i=1}^M$  with the corresponding payments  $\{x_i\}_{i=1}^M$  where  $x_i = b_i \cdot F(b_i)$  and the probability  $\{p_i\}_{i=1}^M$ , where  $\sum p_i = 1$ . The expected payment from this randomized bidding strategy is:

$$p = \sum_{i \in [M]} p_i b_i F(b_i) = \sum_{i \in [M]} p_i x_i .$$

The expected reward  $r$  from this randomized bidding strategy is:

$$r = v \cdot \sum_{i \in [M]} p_i F(b_i) = v \cdot \sum_{i \in [M]} p_i G(x_i) .$$

By the definition of concave envelope, this point  $x = \sum_{i \in [M]} p_i \cdot x_i, y = \sum_{i \in [M]} p_i G(x_i)$  will be inside the convex hull. Therefore, the expected reward from this randomized bidding strategy is at most  $v \cdot G_{\text{concave}}(x)$ .  $\square$

**Lemma A.3.** *Given any  $F$ , defining  $b_0 = \inf\{b \mid F(b) > 0\}$ , there exists a distribution  $F_{\text{conv}} \in \Delta([0, 1])$  such that its allocation-payment curve  $G_{\text{conv}}$  has the following properties:*

- $G_{\text{conv}}(x) = 0$ , for  $x \in [0, b_0 \cdot F(b_0))$  ;
  - $G_{\text{conv}}$  is the concave envelope of the allocation-payment curve  $G$  induced by  $F$ , for  $x \in [b_0 \cdot F(b_0), 1]$  .
- Here,  $x$  represents the input variable of the allocation-payment functions.

*Proof of Theorem A.3.* We prove this lemma statement by explicitly constructing such distribution  $F_{\text{conv}}$ .

Note that the concave envelope of  $G$ , denoted by  $G_{\text{concave}}$ , inherits the monotonic property of  $G$ . Suppose that there is a point  $(x', y')$  on  $G_{\text{concave}}$  where the slope of  $G_{\text{concave}}$  becomes 0. Then for all  $x \geq x'$ , the slope remains 0 due to the convexity of  $G_{\text{concave}}$ . Consequently,  $G_{\text{concave}}$  is constant for  $x \geq x'$ . Since we know  $G_{\text{concave}}(1) = 1$ , it follows that  $G_{\text{concave}}(x) = 1$  for all  $x \geq x'$ .

Hence, considering the definition of  $G_{\text{conv}}$  and the property of the concave envelope, we can split it into three parts:

- For  $x \in [0, b_0 \cdot F(b_0))$ ,  $G_{\text{conv}}(x) = 0$  ;
- For  $x \in [b_0 \cdot F(b_0), x']$ ,  $G_{\text{conv}} = G_{\text{concave}}$  is strictly increasing and thus invertible. We define the invertible function of  $G_{\text{conv}}(\cdot)$  as

$$\beta(y) \triangleq G_{\text{conv}}^{-1}(y), y_0 \triangleq G_{\text{conv}}(b_0 \cdot F(b_0)) = F(b_0) ;$$

- For  $x \in [x', 1]$ ,  $G_{\text{conv}}(x) = 1$ .

The concrete construction process is split into 3 parts based on the value of  $x$ :

**C1: When  $x \in [0, b_0 \cdot F(b_0))$ .** By the definition of  $G_{\text{conv}}(x)$ , we can easily derive that  $F_{\text{conv}}(b) = 0$  and  $b \in [0, b_0)$ ;

**C2: When  $x \in [b_0 \cdot F(b_0), x')$ .** Assume that  $F_{\text{conv}}$  exists, for any  $b_0 \leq b_1 < b_2 < x'$ , if  $F(b_1) = F(b_2)$ , then  $G_{\text{conv}}(b_1 \cdot F_{\text{conv}}(b_1)) = G_{\text{conv}}(b_2 \cdot F_{\text{conv}}(b_2))$  while  $b_1 \cdot F_{\text{conv}}(b_1) \neq b_2 F_{\text{conv}}(b_2)$ , which violates the strict monotonicity of  $G_{\text{conv}}$ . So we can infer that for any  $b_0 \leq b_1 < b_2 < x'$ ,  $F(b_1) \neq F_{\text{conv}}(b_2)$ . This proves that if  $F_{\text{conv}}$  exists, the invertible function of  $F_{\text{conv}}$  exists.

When  $x \in [0, x')$ ,  $G(x) \in [y_0, 1]$ . Since  $\beta(y) = b \cdot y$ , we have

$$b = \frac{\beta(y)}{y} = F_{\text{conv}}^{-1}(y) ,$$

therefore,

$$F_{\text{conv}}^{-1}(y) = \frac{\beta(y)}{y} .$$

Define  $\gamma(y) \triangleq \frac{\beta(y)}{y}$ . Since  $\beta(y)$  is known, we can get  $F_{\text{conv}}^{-1}(y) = \gamma(y)$ . Note that  $\gamma(y)$  must be strictly monotonic in  $y$  in order for  $F_{\text{conv}}$  to exist. Indeed, if  $\gamma$  fails to be consistently monotonic, then we cannot assign a unique  $d$  to each  $y$ , and thus cannot construct a well-defined  $F_{\text{conv}}$ . When  $\gamma$  is strictly monotonic, we can define its inverse  $\gamma^{-1}$ . Consequently, on any interval where  $\beta(y)/y$  is strictly increasing, we obtain

$$F_{\text{conv}}(b) = \gamma^{-1}(b) .$$

Our next step is to prove  $\frac{\beta(y)}{y}$  is strictly increasing. Set  $y_0 < a_1 < a_2 \leq 1$ , then  $k \in (0, 1)$  exists such that  $a_1 = k \cdot a_2 + (1 - k) \cdot y_0$ . Due to the convexity of  $\beta(y)$  and  $\beta(y_0) = 0$

$$\frac{\beta(a_1)}{a_1} = \frac{\beta(k \cdot a_2 + (1 - k) \cdot y_0)}{k \cdot a_2 + (1 - k) \cdot y_0} \frac{k \cdot \beta(a_2) + (1 - k) \cdot \beta(y_0)}{k \cdot a_2 + (1 - k) \cdot y_0} \leq \frac{k \cdot \beta(a_2)}{k \cdot a_2 + (1 - k) \cdot y_0} .$$

Since  $y_0 = F(b_0) > 0$  from the definition of  $b_0$ ,

$$\frac{\beta(a_1)}{a_1} < \frac{k \cdot \beta(a_2)}{k \cdot a_2} = \frac{\beta(a_2)}{a_2} .$$

We find that  $\frac{\beta(y)}{y}$  is strictly increasing, therefore,  $F_{\text{conv}}(b)$  exists, for  $x \in [b_0 \cdot F(b_0), x')$

**C3: When  $x \in [x', 1]$**  Since  $G_{\text{conv}} = 1$ , we have

$$G_{\text{conv}}(x) = F_{\text{conv}}(b) = 1 ,$$

when  $b = b F_{\text{conv}}(b) = x$ .

Combining C1 and C2, we obtain a piecewise definition of  $F_{\text{conv}}(b)$ :

$$F_{\text{conv}}(b) = \begin{cases} 0 , & b \in [0, b_0) , \\ \gamma^{-1}(b) , & b \in [b_0, x') , \\ 1 , & b \in [x', 1] . \end{cases}$$

Here,  $b_0$  is defined in Theorem A.3;  $x'$  is the variable of concave envelope  $G_{\text{concave}}$  such that  $(x', G_{\text{concave}}(x'))$  is the point where the slope of  $G_{\text{concave}}$  becomes 0.

The proof then finishes.  $\square$

**Lemma A.4.** *Fix any value  $v \in [0, 1]$  and any competing bid distribution  $F$ . Let  $F_{\text{conv}}$  be defined as in Theorem A.3. Given any bid  $b \in [0, 1]$ , there exists a distribution over two bids  $b_1, b_2$  with probability  $p_1, p_2$  (respectively), such that the bidder's expected reward (resp. expected payment) when submitting the deterministic bid  $b$  under  $F_{\text{conv}}$  matches the expected reward (resp. expected payment) when submitting the randomized bids  $b_1$  w.p.  $p_1$  and  $b_2$  w.p.  $p_2$  under  $F$ , namely, we have  $vF_{\text{conv}}(b) = p_1 \cdot vF(b_1) + p_2 \cdot vF(b_2)$  (resp.  $bF_{\text{conv}}(b) = p_1 \cdot b_1 F(b_1) + p_2 \cdot b_2 F(b_2)$ ).*

*Proof of Theorem A.4.* According to the definition of  $F_{\text{conv}}$  in Theorem A.3, when  $b \in [0, b_0)$  ( $b_0 = \inf\{b \mid F(b) > 0\}$ ), both  $F_{\text{conv}}(b)$  and  $F(b)$  equal zero, which cause zero payment and reward. Thus, either the (possibly randomized) optimal strategy in  $F$  or the deterministic strategy in  $F_{\text{conv}}$  will not select them as bids. For the proof of this lemma, we can consider only the bid where  $b \in [b_0, 1]$ , that is, the phase in the concave envelope  $G_{\text{concave}}$ .

After we calculate the optimal  $b_{\text{conv}}$  with  $F_{\text{conv}}$  and obtain the corresponding point  $(x_{\text{conv}}, y_{\text{conv}})$  on the concave envelope  $G_{\text{concave}}$ , we can express this point as a convex combination of points on the original curve  $G$ . Specifically, we can find two bids  $x_1, x_2 \in [b_0 \cdot F(b_0), 1]$  and the non-negative coefficient  $p_1$  and  $p_2$  satisfying  $p_1 + p_2 = 1$ , such that

$$(x_{\text{conv}}, y_{\text{conv}}) = (p_1 x_1 + p_2 x_2, p_1 G(x_1) + p_2 G(x_2)) .$$

Since we know  $y_1 = G(x_1) = G_{\text{concave}}(x_1)$  and  $y_2 = G(x_2) = G_{\text{concave}}(x_2)$ , combining with  $x_{\text{conv}} = p_1 x_1 + p_2 x_2$  and  $p_1 + p_2 = 1$ ,  $x_1, x_2$  can be determined.

After obtaining  $(x_1, G(x_1))$  and  $(x_2, G(x_2))$ , we can further determine  $b_1 = \frac{x_1}{G(x_1)}$  with  $F(b_1) = G(x_1)$ , and  $b_2 = \frac{x_2}{G(x_2)}$  with  $F(b_2) = G(x_2)$ . Thus, bids  $b_1$  and  $b_2$  with probability  $p_1$  and  $p_2$  define the optimal randomized strategy on the original competing bid distribution  $F$ , and it achieves the same expected reward and expected payment with the deterministically bidding  $b_{\text{conv}}$  under the  $F_{\text{conv}}$ :

$$vF_{\text{conv}}(b) = v \cdot (p_1 y_1 + p_2 y_2) = p_1 \cdot vF(b_1) + p_2 \cdot vF(b_2) ,$$

and

$$bF_{\text{conv}}(b) = p_1 x_1 + p_2 x_2 = p_1 \cdot b_1 F(b_1) + p_2 \cdot b_2 F(b_2) .$$

$\square$

## A.2 Proofs of Theorem 3.2

*Proof of ROI violation in Theorem 3.2.* We discuss the expected ROI violation by first fixing a value sequence  $V^T$  and then calculating the expected value when considering any value sequence  $V^T$  i.i.d. drawn from  $H^T$ .

For a fixed value sequence  $V^T$ , we have  $b_t = \operatorname{argmax}_{b \in [0, 1]} \left[ \frac{1+\lambda_t}{\lambda_t} \cdot v_t \cdot F_{\text{conv}}(b) - b \cdot F_{\text{conv}}(b) \right]$  at each round  $t$ . Since  $b \in [0, 1]$ , when  $b = 0$ , the function  $\frac{1+\lambda_t}{\lambda_t} \cdot v_t \cdot F_{\text{conv}}(b) - b \cdot F_{\text{conv}}(b) \geq 0$ . Therefore,  $\max \left[ \frac{1+\lambda_t}{\lambda_t} \cdot v_t \cdot F_{\text{conv}}(b) - b \cdot F_{\text{conv}}(b) \right] \geq 0$ . Thus,

$$\frac{1 + \lambda_t}{\lambda_t} \cdot v_t \cdot F_{\text{conv}}(b_t) - b_t \cdot F_{\text{conv}}(b_t) \geq 0 .$$

We further derive

$$b_t \leq \frac{1 + \lambda_t}{\lambda_t} \cdot v_t .$$

To simplify notation, we denote by  $g_t(b_t)$  the difference between the bidder's utility and payment:

$$g_t(b_t) \triangleq v_t \cdot F_{\text{conv}}(b_t) - b_t \cdot F_{\text{conv}}(b_t) . \quad (9)$$

Therefore,  $g_t(b_t)$  satisfies:

$$g_t(b_t) = v_t \cdot F_{\text{conv}}(b_t) - b_t \cdot F_{\text{conv}}(b_t) = (v_t - b_t) \cdot F_{\text{conv}}(b_t) \geq -\frac{v_t}{\lambda_t} \cdot F_{\text{conv}}(b_t) \geq -\frac{1}{\lambda_t} .$$

Thus,

$$\max\left(-1, -\frac{1}{\lambda_t}\right) \leq g_t(b_t) \leq v_t \cdot F_{\text{conv}}(b_t) .$$

Given  $\alpha = \frac{1}{\sqrt{T}}$  and the total ROI violation  $\Delta(\text{Algorithm 1} | (F, V^T)) = -\sum_{t \in [T]} g_t(b_t)$ :

- If  $-\sum_{t \in [T]} g_t(b_t) \leq \log T \sqrt{T}$ , the total ROI violation is bounded by:

$$\Delta(\text{Algorithm 1} | (F, V^T)) = -\sum_{t \in [T]} g_t(b_t) \leq \log T \sqrt{T} .$$

- If  $-\sum_{t \in [T]} g_t(b_t) > \log T \sqrt{T}$ , we can find  $T'$  which is the last time that  $\sum_{t \in [T]} g_t(b_t) \leq \log T \sqrt{T}$ , so we know for any  $t > T' + 1$ , the dual variable  $\lambda_t$  must be larger than  $T$  since

$$\lambda_t = \exp\left[-\alpha \sum_{t'=1}^{t-1} g'_t(b'_t)\right] > \exp[\alpha \sqrt{T}] = T .$$

Thus,

$$-g_t(b_t) \leq \frac{1}{\lambda_t} \leq \frac{1}{T} .$$

The ROI violation is bounded by

$$\begin{aligned} \Delta(\text{Algorithm 1} | (F, V^T)) &= -\sum_{t \in [T]} g_t(b_t) \\ &= -\sum_{t \in [T']} g_t(b_t) + -\sum_{t > T'} g_t(b_t) \\ &\leq \log T \sqrt{T} + 2 \\ &\leq 2 \log T \sqrt{T} . \end{aligned}$$

Thus, for  $V^T$  i.i.d. drawn from  $H^T$ , the expected ROI violation of the algorithm under the distribution  $H$  satisfies

$$\Delta(\text{Algorithm 1} | (F, H)) = \mathbb{E}_{V^T \sim H^T} [\Delta(\text{Algorithm 1} | (F, V^T))] \leq 2 \log T \sqrt{T} .$$

□

*Proof of expected regret in Theorem 3.2.* We split the proof into two main steps:

**Step 1: To prove the regret is bounded by  $\mathbb{E}_{V^T \sim H^T} \left[ \sum_{t \in [T]} \lambda_t \cdot g_t(b_t) \right]$ .** The reward function  $r_t(b) = v_t \cdot F(b)$ . For  $\lambda \geq 0$ , define a function of  $\lambda$  as:

$$f_t^*(\lambda) \triangleq \max_{b \in [0,1]} [r_t(b) + \lambda \cdot g_t(b)] ,$$

where  $g_t(b)$  is defined in Eqn. 9. Then, we define

$$\bar{D}(\lambda | (F, H)) \triangleq \mathbb{E}_{V^T \sim H^T} [f_t^*(\lambda)] .$$

We first get the upper bound of the optimal randomized strategy  $\text{OPT}$ . Due to the weak duality, we know that the primal question will be less or equal to the dual problem. Since  $v_t$  is i.i.d generated from the distribution  $H$ , the optimal strategy of each round is independent. We can prove that the total reward of  $T$  rounds for the optimal strategy will be bounded by:

$$\text{Reward}(\text{OPT} | (F, H)) \leq T \cdot \min_{\lambda \geq 0} \bar{D}(\lambda | (F, H)) ,$$

where  $\text{OPT}$  is the corresponding optimal bidding algorithm in hindsight. Second, we prove the lower bound of Algorithm 1. when we get the dual variable  $\lambda_t$  at round  $t$ , the lower bound of Algorithm 1's reward is

$$\text{Reward}(\text{Algorithm 1} | (F, H)) = T \cdot \bar{D}(\bar{\lambda}_T | (F, H)) - \mathbb{E}_{V^T \sim H^T} \left[ \sum_{t \in [T]} \lambda_t \cdot g_t(b_t) \right] ,$$

where  $\bar{\lambda}_T = \frac{1}{T} \sum_{t \in [T]} \lambda_t$ .

To prove this, at round  $t$ , from the optimization algorithm, we can get that

$$b_t = \text{argmax}[v_t \cdot F(b) + \lambda_t \cdot g_t(b)] .$$

Therefore, the reward in round  $t$  satisfies

$$\begin{aligned} r_t(b_t) &= v_t \cdot F(b_t) \\ &= \max [v_t \cdot F(b) + \lambda_t \cdot g_t(b)] - \lambda_t \cdot g_t(b_t) \\ &= f_t^*(\lambda_t) - \lambda_t \cdot g_t(b_t) . \end{aligned}$$

Taking expectation on both sides given the output until iteration  $t-1$ , denote the condition as  $\sigma_{t-1}$ :

$$\mathbb{E}[r_t(b_t) | \sigma_{t-1}] = \mathbb{E}[f_t^*(\lambda_t) | \sigma_{t-1}] - \mathbb{E}[\lambda_t \cdot g_t(b_t) | \sigma_{t-1}] .$$

Because  $\sigma_{t-1}$  decides  $\lambda_t$ ,  $\mathbb{E}[f_t^*(\lambda_t) | \sigma_{t-1}] = \mathbb{E}_{V^{t-1} \sim H^{t-1}} [f_t^*(\lambda_t)]$ . We have

$$\mathbb{E}[r_t(b_t) | \sigma_{t-1}] = \bar{D}(\lambda_t | (F, H)) - \mathbb{E}[\lambda_t \cdot g_t(b_t) | \sigma_{t-1}] .$$

The sum of the rewards over  $T$  rounds,

$$\begin{aligned} \text{Reward}(\text{Algorithm 1} | (F, H)) &= \mathbb{E}_{V^T \sim H^T} \left[ \sum_{t \in [T]} r_t(b_t) \right] \\ &= \sum_{t \in [T]} \bar{D}(\lambda_t | (F, H)) - \mathbb{E}_{V^T \sim H^T} \left[ \sum_{t \in [T]} \lambda_t \cdot g_t(b_t) \right] . \end{aligned}$$

By the Definition of regret in Eqn. (1), we have:

$$\begin{aligned}
\text{Regret}(\text{Algorithm 1} \mid (F, H)) &= \mathbb{E}_{V^T \sim H^T} \left[ \text{Reward(OPT} \mid (F, V^T)) - \text{Reward(Algorithm 1} \mid (F, V^T) \right] \\
&= \mathbb{E}_{V^T \sim H^T} \left[ \text{Reward(OPT} \mid (F, V^T) \right] - \mathbb{E}_{V^T \sim H^T} \left[ \text{Reward(Algorithm 1} \mid (F, V^T) \right] \\
&\leq T \cdot \min_{\lambda \geq 0} \bar{D}(\lambda \mid (F, H)) - \sum_{t \in [T]} \bar{D}(\lambda_t \mid (F, H)) + \mathbb{E}_{V^T \sim H^T} \left[ \sum_{t \in [T]} \lambda_t \cdot g_t(b_t) \right] \\
&\leq \mathbb{E}_{V^T \sim H^T} \left[ \sum_{t \in [T]} \lambda_t \cdot g_t(b_t) \right] .
\end{aligned}$$

**Step 2: To prove, for any sequence of input value  $V^T$ , the sequences of  $\{b_t\}_{t=1}^T$  and  $\{\lambda_t\}_{t=1}^T$  generated by Algorithm 1 satisfies**

$$\sum_{t \in [T]} \lambda_t \cdot g_t(b_t) = O(\sqrt{T}) . \quad (10)$$

Since  $\lambda_{t+1} = \lambda_t \cdot \exp[-\alpha \cdot g_t(b_t)]$ ,  $\alpha g_t$  can be represented as:

$$\alpha g_t(b_t) = \log(\lambda_t / \lambda_{t+1}) .$$

In the mirror descent framework, we have  $V_h(y, x) = h(y) - h(x) - h'(x) \cdot (y - x)$  (Bregman divergence) and  $h(u) = u \log u - u$  (generalized negative entropy); thus,  $V_h(y, x)$  can be rewritten as:

$$V_h(y, x) = y \log(y/x) - y + x .$$

Therefore, we can rewrite  $\alpha g_t(b_t) \cdot \lambda_t$  as:

$$\begin{aligned}
\alpha g_t(b_t) \cdot \lambda_t &= \alpha g_t(b_t) \cdot (\lambda_t - \lambda_{t+1}) + \alpha g_t(b_t) \cdot \lambda_{t+1} \\
&= \alpha g_t(b_t) \cdot (\lambda_t - \lambda_{t+1}) + \log(\lambda_t / \lambda_{t+1}) \cdot \lambda_{t+1} \\
&= (\alpha g_t(b_t) + 1) \cdot (\lambda_t - \lambda_{t+1}) - V_h(\lambda_{t+1}, \lambda_t) .
\end{aligned}$$

Due to the local strong convexity of  $V_h$ , for any  $x, y > 0$ ,

$$V_h(y, x) = y \log(y/x) - y + x \geq \frac{1}{2 \max(x, y)} \cdot (y - x)^2 ,$$

so  $\alpha g_t(b_t) \cdot \lambda_t$  satisfies:

$$\alpha g_t(b_t) \cdot \lambda_t \leq (\alpha g_t(b_t) + 1) \cdot (\lambda_t - \lambda_{t+1}) - \frac{(\lambda_{t+1} - \lambda_t)^2}{2 \max(\lambda_t, \lambda_{t+1})} .$$

Based on Cauchy-Schwarz inequality,

$$\left( \frac{\lambda_t - \lambda_{t+1}}{\sqrt{2} \max(\lambda_t, \lambda_{t+1})} \right)^2 + \left( \frac{\alpha g_t(b_t)}{\sqrt{2}} \right)^2 \geq \alpha g_t(b_t) \cdot \frac{\lambda_t - \lambda_{t+1}}{\max(\lambda_t, \lambda_{t+1})} .$$

Therefore,

$$\frac{(\lambda_t - \lambda_{t+1})^2}{2 \max(\lambda_t, \lambda_{t+1})} + \frac{1}{2} \alpha^2 g_t^2(b_t) \cdot \max(\lambda_t, \lambda_{t+1}) \geq \alpha g_t(b_t) \cdot (\lambda_t - \lambda_{t+1}) .$$

$\alpha g_t(b_t) \cdot \lambda_t$  satisfies:

$$\alpha g_t(b_t) \cdot \lambda_t \leq (\lambda_t - \lambda_{t+1}) + \frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1}). \quad (11)$$

We know that  $g_t(b_t)$  is bounded by  $\max(-1, -1/\lambda_t) \leq g_t(b_t) \leq v_t \cdot x_t(b_t)$ . We now bound the term  $\frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1})$  in a case-wise manner.

**C1: When we assume**  $g_t(b_t) \geq 0$ .  $g_t(b_t) \leq 1$  with  $\alpha = \frac{1}{\sqrt{T}}$  imply  $\alpha g_t(b_t) \leq 1$ , so  $\exp(-\alpha g_t(b_t)) \leq 1 - \frac{\alpha g_t(b_t)}{2}$

$$\lambda_{t+1} = \lambda_t \exp[-\alpha g_t(b_t)] \leq \lambda_t \cdot (1 - \alpha g_t(b_t)/2) = \lambda_t - \frac{1}{2} \alpha \lambda_t \cdot g_t(b_t).$$

We can infer that  $\lambda_t \geq \lambda_{t+1}$ . Combine with the bound of  $g_t(b_t)$ :  $0 \leq g_t(b_t) \leq 1$  to obtain:

$$\frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1}) = \frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \lambda_t \leq \alpha \frac{\alpha g_t(b_t) \lambda}{2} \leq \alpha(\lambda_t - \lambda_{t+1}).$$

**C2: When we assume**  $g_t(b_t) < 0$ .  $g_t(b_t) \geq \max(-1, -\frac{1}{\lambda_t})$  and  $\lambda_t \leq \lambda_{t+1}$

$$\frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1}) = \frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \lambda_{t+1} \leq \alpha^2 \cdot \frac{1}{\lambda_t} \cdot \lambda_{t+1}.$$

Since  $g_t(b_t) \geq -1$  with  $\alpha = \frac{1}{\sqrt{T}}$ , it implies  $\alpha g_t(b_t) \geq -1$ . Therefore,

$$\lambda_{t+1} = \lambda_t \exp[-\alpha g_t(b_t)] \leq e \lambda_t.$$

We have

$$\frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1}) \leq 2\alpha^2.$$

When  $\alpha g_t(b_t) \in [-1, 0]$ ,  $\exp[-\alpha g_t(b_t)] \leq (1 - 2\alpha g_t(b_t))$ . Therefore, the relationship between  $\lambda_{t+1}$  and  $\lambda_t$  satisfies

$$\lambda_{t+1} = \lambda_t \exp[-\alpha g_t(b_t)] \leq \lambda_t(1 - 2\alpha g_t(b_t)).$$

Applying  $g_t(b_t) \in [-\frac{1}{\lambda_t}, 0]$ ,

$$\lambda_{t+1} - \lambda_t \leq -2\lambda_t \cdot \alpha g_t(b_t) \leq 2\alpha.$$

Thus, Summarizing different cases,  $\alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1})$  satisfies:

$$\frac{1}{2} \alpha^2 g_t(b_t)^2 \cdot \max(\lambda_t, \lambda_{t+1}) \leq \alpha(\lambda_t - \lambda_{t+1}) + 2\alpha^2.$$

Since  $\alpha = \frac{1}{\sqrt{T}} \geq 1$  when  $T \geq 1$ , combining with Eqn. (11), we have

$$\alpha g_t(b_t) \cdot \lambda_t \leq 2(\lambda_t - \lambda_{t+1}) + 2\alpha^2.$$

Sum over  $T$  rounds:

$$\begin{aligned} \sum_{t \in [T]} g_t(b_t) \cdot \lambda_t &\leq \frac{1}{\alpha} \sum_{t \in [T]} (2(\lambda_t - \lambda_{t+1}) + 2\alpha^2) \\ &\leq \frac{1}{\alpha} (4 - 2\lambda_{T+1}) \\ &= O(\sqrt{T}). \end{aligned}$$

The Eqn. (10) has been claimed.  $\square$

## B Missing Proofs in Section 4

**Lemma B.1.** *For an unknown true CDF  $F$ , we construct an estimated CDF  $\hat{F}_m$  as in Eqn. (4) and define an optimistic CDF  $\bar{F}_m$  as in Eqn. (5). Consequently, the distance between  $\bar{F}_m$  and  $F$  can be bounded by  $2\varepsilon_m$  and  $\bar{F}_m(b) \geq F(b)$  for  $b \in [0, 1]$  with probability  $1 - O(1/N_m)$ , when  $\varepsilon_m$  satisfies  $\varepsilon_m = \Theta\left(\frac{\log N_m}{\sqrt{N_m}}\right)$  ( $N_m$  is the number of samples we used for estimation).*

*Proof of Theorem B.1.* Based on the definition of the estimated CDF  $\hat{F}$ , Dvoretzky–Kiefer–Wolfowitz inequality bounds the probability that the CDF  $\hat{F}(b)$  differs from  $F(b)$  by more than a given  $\varepsilon_m > 0$ :

$$\Pr\left[\sup_{b \in [0, 1]} |\hat{F}_N(b) - F(b)| > \varepsilon_m\right] \leq \exp(-2N_m\varepsilon_m^2).$$

To ensure  $\Pr\left[\sup_{b \in [0, 1]} |\hat{F}_m(b) - F(b)| > \varepsilon_m\right]$  is sufficiently small (i.e. approach zero as  $N_m$  increases),  $2N_m\varepsilon_m^2$  must tend to infinity. Therefore, when  $\varepsilon_m$  satisfies

$$\varepsilon_m = \Theta\left(\frac{\log N_m}{\sqrt{N_m}}\right),$$

the certainty of the closeness increases as  $N_m$  increases. Because  $F$  is within  $\varepsilon_m$  of  $\hat{F}_m$  and  $\bar{F}_m$  is defined as  $(\hat{F}_m(b) + \varepsilon_m) \wedge 1$ , it follows that

$$\sup_b |F(b) - \bar{F}_m(b)| \leq 2\varepsilon,$$

and

$$F(b) \leq \hat{F}_m(b) + \varepsilon = \bar{F}_m(b),$$

for all  $b$ . □

### B.1 Analysis of Algorithm 1 with Estimated CDF

This section focuses on the analysis of Algorithm 1 with optimistic empirical CDF  $\bar{F}$  relative to a true CDF  $F$ . As a result, Theorem B.2 proved that near-optimal performance (Regret and ROI violation) is still guaranteed with the support of Theorem B.3 and Theorem B.4.

**Lemma B.2.** *Fix any value sequence  $V^T$ . Given the two CDFs,  $F$  and  $F^\dagger$ , such that  $F^\dagger(b) \geq F(b)$  for any  $b \in [0, 1]$  and  $\sup_b |F(b) - F^\dagger(b)| \leq 2\varepsilon$ . Run Algorithm 1 with the input CDF  $F^\dagger$  (denoted by  $\text{ALG1}(F^\dagger)$ ) with  $F$  as environment for  $T$  rounds. The Regret satisfies:*

$$\text{Regret}(\text{ALG1}(F^\dagger) \mid (F, V^T)) \leq O(\sqrt{T}) + T \cdot 2\varepsilon,$$

and the ROI violation satisfies:

$$\Delta(\text{ALG1}(F^\dagger) \mid (F, V^T)) \leq 2\sqrt{T} \log T + T \cdot 2\varepsilon,$$

*Proof of Theorem B.2.* From Theorem B.3 and Theorem B.4, we have

$$\text{Reward}(\text{ALG1}(F^\dagger) \mid (F, V^T)) \geq \text{Reward}(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) - T \cdot 2\varepsilon$$

and

$$\text{Reward}(\text{OPT}(F^\dagger) | (F^\dagger, V^T)) \geq \text{Reward}(\text{OPT}(F) | (F, V^T)) .$$

Here, we denote output of the optimal bidding algorithm in hindsight with  $F^\dagger$  as input by  $\text{OPT}(F^\dagger)$ . Similarly, we denote the optimal bidding algorithm in hindsight with  $F$  as input by  $\text{OPT}(F)$ .

Since the Regret of Algorithm 1 with known true CDF is  $O(\sqrt{T})$ , we can infer that

$$\text{Reward}(\text{ALG1}(F^\dagger) | (F^\dagger, V^T)) = \text{Reward}(\text{OPT}(F^\dagger) | (F^\dagger, V^T)) - O(\sqrt{T}) .$$

Combining above, we obtain

$$\begin{aligned} \text{Reward}(\text{ALG1}(F^\dagger) | (F, V^T)) &\geq \text{Reward}(\text{ALG1}(F^\dagger) | (F^\dagger, V^T)) - T \cdot 2\varepsilon \\ &= \text{Reward}(\text{OPT}(F^\dagger) | (F^\dagger, V^T)) - T \cdot 2\varepsilon - O(\sqrt{T}) \\ &\geq \text{Reward}(\text{OPT}(F) | (F, V^T)) - T \cdot 2\varepsilon - O(\sqrt{T}) . \end{aligned}$$

Therefore, the Regret satisfies

$$\begin{aligned} \text{Regret}(\text{ALG1}(F^\dagger) | (F, V^T)) &= \text{Reward}(\text{OPT}(F) | (F, V^T)) - \text{Reward}(\text{ALG1}(F^\dagger) | (F, V^T)) \\ &= O(\sqrt{T}) + T \cdot 2\varepsilon . \end{aligned}$$

For ROI violation, combining

$$\Delta(\text{ALG1}(F^\dagger) | (F, V^T)) \leq \Delta(\text{ALG1}(F^\dagger) | (F^\dagger, V^T)) + T \cdot 2\varepsilon$$

and

$$\Delta(\text{ALG1}(F^\dagger) | (F^\dagger, V^T)) \leq 2\sqrt{T} \log T ,$$

We have

$$\Delta(\text{ALG1}(F^\dagger) | (F, V^T)) \leq 2\sqrt{T} \log T + T \cdot 2\varepsilon . \quad \square$$

**Lemma B.3.** *Fix any value sequence  $V^T$ . Given  $F^\dagger$  and  $F$  are CDFs of maximum competing bid, when  $F^\dagger(b) \geq F(b)$  for any  $b \in [0, 1]$ , we claim that this implies the Reward of the optimal bidding algorithm in hindsight with CDF input  $F^\dagger$  (denoted by  $\text{OPT}(F^\dagger)$ ) and  $F$  (denoted by  $\text{OPT}(F)$ ) satisfies:*

$$\text{Reward}(\text{OPT}(F^\dagger) | (F^\dagger, V^T)) \geq \text{Reward}(\text{OPT}(F) | (F, V^T)) .$$

*Proof of Theorem B.3.* Construct the concave envelope  $G_{\text{conv}}$  and  $G_{\text{conv}}^\dagger$ :

- Given  $F(b)$ , define a curve  $G$  as

$$y = G(x), y = F(b), x = b \cdot F(b) .$$

Let  $G_{\text{conv}}$  denote the concave envelope of  $G$ .

- Similarly, for  $F^\dagger(b)$ , define a curve  $G^\dagger$ , where

$$y = G^\dagger(x), y = F^\dagger(b), x = b \cdot F^\dagger(b) ,$$

and denote its concave envelope as  $G_{\text{conv}}^\dagger$ .

Due to the property of the cumulative distribution function,  $F(b)$  and  $F^\dagger(b)$  both have definition at  $b = 0$  and  $F^\dagger(0) \geq F(0) \geq 0$ .

Fix a bit  $b$ , there exists a point  $(x, y)$  the curve  $G$ :

$$x = b \cdot F(b), y = F(b) ,$$

and a corresponding point  $(x^\dagger, y^\dagger)$  on the curve  $G^\dagger$ :

$$x^\dagger = b \cdot F^\dagger(b), y^\dagger = F^\dagger(b) .$$

Notice that both points lie on the line:  $y = \frac{1}{b}x$ . Because  $F^\dagger(b) \geq F(b)$ ,  $\frac{F(b)}{F^\dagger(b)} \in [0, 1]$ . Observe that the point  $(x, y)$  is a convex combination of  $(x^\dagger, y^\dagger)$  and  $(0, 0)$ :

$$\begin{aligned} x &= \frac{F(b)}{F^\dagger(b)} \cdot x^\dagger + \left(1 - \frac{F(b)}{F^\dagger(b)}\right) \cdot 0 , \\ y &= \frac{F(b)}{F^\dagger(b)} \cdot y^\dagger + \left(1 - \frac{F(b)}{F^\dagger(b)}\right) \cdot 0 . \end{aligned}$$

This shows that every point on the curve  $G$  is contained in the convex hull of  $G^\dagger \cup \{(0, 0)\}$ . Due to the property of CDF,  $F(b)$  and  $F^\dagger(b)$  both have definition at  $b = 0$  and  $F^\dagger(0) \geq F(0) \geq 0$ .  $(0, 0)$  is in both the convex hull of  $G$  and the convex hull of  $G^\dagger$ . Thus, the convex hull of  $G$  is a subset of the convex hull of  $G^\dagger$ .

By the definition of concave envelope, the lowest-valued concave function that overestimates or equals the original function over that set, it follows

$$G_{\text{conv}} \leq G_{\text{conv}}^\dagger ,$$

pointwise.

Thus, for a fixed value sequence  $V_T$ , given CDF  $F$ , an optimal bid sequence  $\{b_t\}_{t=1}^T$  will be generated. For each  $b_t$  at round  $t$ , we could find  $b_t^\dagger$  under  $F^\dagger$  with the same payment ( $b_t F(b_t) = b_t^\dagger F(b_t^\dagger)$ ) but a potentially higher value ( $v_t \cdot F^\dagger(b_t^\dagger) \geq v_t \cdot F(b_t)$ ). Therefore, the reward satisfies

$$\begin{aligned} \text{Reward}(\text{OPT}(F) \mid (F, V^T)) &= \sum_{t \in [T]} v_t \cdot F(b_t) \\ &\leq \sum_{t \in [T]} v_t \cdot F^\dagger(b_t^\dagger) \\ &= \text{Reward}(\text{OPT}(F^\dagger) \mid (F^\dagger, V^T)) . \end{aligned} \quad \square$$

**Lemma B.4.** *Fix a value sequence  $V^T$ . Given the two CDFs,  $F$  and  $F^\dagger$ , such that  $F^\dagger(b) \geq F(b)$  and  $\sup_b |F(b) - F^\dagger(b)| \leq 2\varepsilon$ . Run Algorithm 1 with  $F^\dagger$  as its input, denoted by  $\text{ALG1}(F^\dagger)$ . The reward satisfies*

$$\text{Reward}(\text{ALG1}(F^\dagger) \mid (F, V^T)) \geq \text{Reward}(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) - T \cdot 2\varepsilon ,$$

and the ROI violation cannot be amplified by more than  $T \cdot 2\varepsilon$ :

$$\Delta(\text{ALG1}(F^\dagger) \mid (F, V^T)) \leq \Delta(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) + T \cdot 2\varepsilon .$$

*Proof of Theorem B.4.* At round  $t$ , given the input value  $v_t$ , the bidding algorithm  $\text{ALG1}(F^\dagger)$  generates  $b_{\text{conv}}^\dagger$  and the randomized strategy  $b_1^\dagger$  with probability  $p_1$  and  $b_2^\dagger$  with probability  $p_2$ .  $p_1 + p_2 = 1$  and  $p_1 \geq 0, p_2 \geq 0$ .

Under  $F^\dagger$ , it generates expected reward

$$r_t^\dagger = v_t \cdot [p_1 \cdot F^\dagger(b_1^\dagger) + p_2 \cdot F^\dagger(b_2^\dagger)] ,$$

and expected payment

$$p_t^\dagger = p_1 \cdot b_1^\dagger \cdot F^\dagger(b_1^\dagger) + p_2 \cdot b_2^\dagger \cdot F^\dagger(b_2^\dagger) .$$

When the distribution is  $F$ , it generates expected reward

$$r_t = v_t \cdot [p_1 \cdot F(b_1^\dagger) + p_2 \cdot F(b_2^\dagger)] ,$$

and expected payment

$$p_t = p_1 \cdot b_1^\dagger \cdot F(b_1^\dagger) + p_2 \cdot b_2^\dagger \cdot F(b_2^\dagger) .$$

Because  $\sup_b |F(b) - F^\dagger(b)| \leq 2\varepsilon$  and  $F(b) \leq F^\dagger(b)$ , The difference in value will be  $F(b) \geq F^\dagger(b) - 2\varepsilon$ .

The difference of expected reward will be

$$\begin{aligned} r_t - r_t^\dagger &= v_t \cdot \{p_1 \cdot [F(b_1^\dagger) - F^\dagger(b_1^\dagger)] + p_2 \cdot [F(b_2^\dagger) - F^\dagger(b_2^\dagger)]\} \\ &\geq v_t \cdot [p_1 \cdot (-2\varepsilon) + p_2 \cdot (-2\varepsilon)] \\ &\geq -2\varepsilon . \end{aligned}$$

The difference of expected payment will be

$$p_t - p_t^\dagger = p_1 \cdot b_1^\dagger \cdot [F(b_1^\dagger) - F^\dagger(b_1^\dagger)] + p_2 \cdot b_2^\dagger \cdot [F(b_2^\dagger) - F^\dagger(b_2^\dagger)] \leq 0 .$$

Sum of  $T$  rounds, the difference of expected reward satisfies:

$$\sum_{t \in [T]} r_t - \sum_{t \in [T]} r_t^\dagger \geq -T \cdot 2\varepsilon ,$$

and the difference of expected payment satisfies:

$$\sum_{t \in [T]} p_t - \sum_{t \in [T]} p_t^\dagger \leq 0 .$$

Therefore, the over all reward satisfies:

$$\begin{aligned} \text{Reward}(\text{ALG1}(F^\dagger) \mid (F, V^T)) &= \sum_{t \in [T]} r_t \\ &\geq \sum_{t \in [T]} r_t^\dagger - T \cdot 2\varepsilon \\ &= \text{Reward}(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) - T \cdot 2\varepsilon . \end{aligned}$$

The ROI violation satisfies:

$$\begin{aligned} \Delta(\text{ALG1}(F^\dagger) \mid (F, V^T)) &= \text{Payment}(\text{ALG1}(F^\dagger) \mid (F, V^T)) - \text{Reward}(\text{ALG1}(F^\dagger) \mid (F, V^T)) \\ &\leq \text{Payment}(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) - \text{Reward}(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) + T \cdot 2\varepsilon \\ &= \Delta(\text{ALG1}(F^\dagger) \mid (F^\dagger, V^T)) + T \cdot 2\varepsilon . \end{aligned} \quad \square$$

## B.2 Analysis of Regret and ROI Violation of Algorithm 2

In this section, we analyze the performance of Algorithm 2 in terms of regret and ROI violation with  $\varepsilon_m = \frac{\log N_m}{\sqrt{N_m}}$  for each stage  $m$ .

**Lemma B.5.** *Consider an unknown CDF  $F$  of maximum competing bid. Given a fixed value sequence  $V^T$ , run Algorithm 2 for  $T$  rounds. With  $\varepsilon_m = \frac{\log N_m}{\sqrt{N_m}}$  for each stage  $m$ , the reward satisfies*

$$\text{Regret}(\text{Algorithm 2} | (F, V^T)) = \tilde{O}(\sqrt{T}) . \quad (12)$$

*Proof of Theorem B.5.* For each stage  $m$ , we run Algorithm 1 with the optimistic CDF  $\bar{F}_m$  for  $T_m$  rounds, where  $N_m = 2^{m-1} - 1$  and  $T_m = 2^{m-1}$ .

Based on Theorem B.2 and  $\varepsilon_m = \frac{\log N_m}{\sqrt{N_m}}$ , we have

$$\begin{aligned} \text{Reward}(\text{ALG1}(\bar{F}_m) | (F, V_m)) &\geq \text{Reward}(\text{OPT} | (F, V_m)) - T_m \cdot 2\varepsilon_m \\ &= \text{Reward}(\text{OPT} | (F, V_m)) - \tilde{O}(\sqrt{T_m}) , \end{aligned} \quad (13)$$

where  $V_m$  is the sub-sequence of values at stage  $m$  including  $T_m$  elements, and  $\text{ALG1}(\bar{F}_m)$  the bids generated by feeding Algorithm 1 with a distribution  $\bar{F}_m$ .

The total reward across  $M$  stages is

$$\begin{aligned} \text{Reward}(\text{Algorithm 2} | (F, V^T)) &= \sum_{m \in [M]} \text{Reward}(\text{ALG1}(\bar{F}_m) | (F, V_m)) \\ &\geq \text{Reward}(\text{OPT}(F) | (F, V^T)) - \sum_{m \in [M]} \sqrt{2^{m-1}} \\ &= \text{Reward}(\text{OPT}(F) | (F, V^T)) - \tilde{O}(\sqrt{T}) . \end{aligned} \quad \square$$

**Lemma B.6.** *Consider an unknown CDF  $F$  of maximum competing bid. Given a fixed value sequence  $V^T$ , run Algorithm 2 for  $T$  rounds. With  $\varepsilon_m = \frac{\log N_m}{\sqrt{N_m}}$  for each stage  $m$ , the ROI Violation satisfies:*

$$\Delta(\text{Algorithm 2} | (F, V^T)) = \tilde{O}(\sqrt{T}) . \quad (14)$$

*Proof of Theorem B.6.* For each stage  $m$ , we run Algorithm 1 with the optimistic CDF  $\bar{F}_m$  for  $T_m = 2^{m-1}$  rounds, where  $N_m = 2^{m-1} - 1$  is the number of samples collected in the preceding stages.

Based on Theorem B.2 and  $\varepsilon_m = \frac{\log N_m}{\sqrt{N_m}}$ , The ROI violation is bounded by

$$\begin{aligned} \Delta(\text{ALG1}(\bar{F}_m) | (F, V_m)) &\leq 2\sqrt{T_m} \log T_m + T_m \cdot 2\varepsilon_m \\ &\leq 2\sqrt{T_m} \cdot (\log T_m + 1) . \end{aligned} \quad (15)$$

The total ROI violation across  $M = \lceil \log_2(T+1) \rceil$  stages is bounded by

$$\begin{aligned} \Delta(\text{Algorithm 2} | (F, V^T)) &= \sum_{m \in [M]} \Delta(\text{ALG1}(\bar{F}_m) | (F, V_m)) \\ &\leq \sum_{m \in [M]} 2\sqrt{2^{m-1}} (\log(2^{m-1}) + 1) \\ &\leq \sum_{m \in [M]} 2 \log 2 \cdot (m-1) \cdot \sqrt{2^{m-1}} + \sum_{m \in [M]} 2\sqrt{2^{m-1}} \\ &= O(M \cdot \sqrt{2^M}) \\ &= O(\sqrt{T} \log T) . \end{aligned} \quad \square$$

## C Missing Proofs in Section 5

### C.1 Construct the Empirical Distribution and its Optimistic Distribution

**Lemma C.1.** *Given the grid points  $\{c_k\}_{k \in [K-1] \cup \{0\}}$  where we bid every point consecutively for  $M$  time rounds, the probability that there exists one grid bid  $c_k$  for which the absolute difference between the estimated CDF  $\widehat{F}(c_k)$  and the true CDF  $F(c_k)$  exceeds  $\varepsilon > 0$  is bounded by*

$$\Pr\left[\exists k \in [K] : |\widehat{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq K \cdot \exp(-2M\varepsilon^2).$$

*Proof of Theorem C.1.* We define  $\widetilde{F}(c_k) \triangleq \frac{1}{M} \cdot \sum_{t \in [k \cdot M + 1 : (k+1) \cdot M]} \mathbf{1}\{b_t \geq d_t\}$ , for notation simplicity. By Chernoff bound, for each grid point  $c_k$  and a given  $\varepsilon > 0$ , we have

$$\Pr\left[|\widetilde{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq \exp(-2M\varepsilon^2).$$

Now consider  $|\widehat{F}(c_k) - F(c_k)|$  where  $i \leq k$ . Since  $\widetilde{F}(c_i) \leq \widehat{F}(c_k)$  and  $F(c_i) \leq F(c_k)$  for all  $i \leq k$ , we have:

$$\begin{aligned} \widehat{F}(c_k) - F(c_k) &= \max_{i \in [k] \cup \{0\}} \{\widetilde{F}(c_k)\} - F(c_k) \\ &= \max_{i \in [k] \cup \{0\}} \{\widetilde{F}(c_i) - F(c_i) + F(c_i)\} - F(c_k) \\ &\leq \max_{i \in [k] \cup \{0\}} \{\widetilde{F}(c_i) - F(c_i)\} + \max_{i \in [k] \cup \{0\}} \{F(c_i)\} - F(c_k) \\ &\leq \max_{i \in [k] \cup \{0\}} \{\widetilde{F}(c_i) - F(c_i)\}. \end{aligned}$$

Similarly,

$$\begin{aligned} F(c_k) - \widehat{F}(c_k) &= F(c_i) - \max_{i \in [k] \cup \{0\}} \{\widetilde{F}(c_i)\} \\ &\leq F(c_k) - \widetilde{F}(c_k) \\ &\leq |F(c_k) - \widetilde{F}(c_k)|. \end{aligned}$$

Therefore,  $|\widehat{F}(c_k) - F(c_k)| \leq \max_{i \in [k] \cup \{0\}} |\widetilde{F}(c_i) - F(c_i)|$ . If  $|\widetilde{F}(c_i) - F(c_i)| < \varepsilon$  for all  $i \leq k$ , then  $|\widehat{F}(c_k) - F(c_k)| < \varepsilon$ .

Thus, the event  $|\widehat{F}(c_k) - F(c_k)| \geq \varepsilon$  implies that there exists at least one  $i \leq k$  such that  $|\widetilde{F}(c_i) - F(c_i)| \geq \varepsilon$ .  $\widetilde{F}(c_i)$  and  $F(c_k)$  is also bounded by

$$\Pr\left[|\widehat{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq \exp(-2M\varepsilon^2).$$

By the union bound,

$$\Pr\left[\exists i \in [K] : |\widehat{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq \sum_i \Pr\left[|\widehat{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq K \cdot \exp(-2M\varepsilon^2). \quad \square$$

**Lemma C.2.** *Given the grid points  $\{c_k\}_{k \in [K] \cup \{0\}}$ , and an empirical CDF  $\bar{F}$  constructed in such that  $\bar{F}(c_k) \geq F(c_k)$  for all  $k \in [K] \cup \{0\}$ ), if we construct another CDF  $\bar{F}^c$  as:*

$$\bar{F}^c(b) = \bar{F}((b + 1/K) \wedge 1),$$

then  $\bar{F}^c(b) \geq F(b)$ ,  $b \in [0, 1]$ .

*Proof of Theorem C.2.* Given that  $\bar{F}(c_k) \geq F(c_k)$  for all  $k \in \{0, 1, \dots, K\}$ . Consider an interval  $[c_k, c_{k+1})$ . For any  $b \in [c_k, c_{k+1})$ , by the definition of  $\bar{F}^c(b)$ , we have

$$\bar{F}^c(b) = \bar{F}(\min(b + 1/K, 1)) .$$

Since  $b \in [c_k, c_{k+1})$ , we have  $b + 1/K \in [c_{k+1}, c_{k+2})$ . Since  $\bar{F}$  is a step function constant in the range  $[c_{k+1}, c_{k+2})$ , we have

$$\bar{F}(\min(b + 1/K, 1)) = \bar{F}(\min(c_{k+1}, 1)) .$$

Both  $\bar{F}$  and  $F$  are non-decreasing functions, we have

$$\bar{F}(\min(c_{k+1}, 1)) \geq F(c_{k+1}) \geq F(b) .$$

Combining these inequalities, we get:

$$\bar{F}^c(b) \geq F(b), b \in [c_k, c_{k+1}) .$$

This holds for all intervals defined by the grid points. Therefore,

$$\bar{F}^c(b) \geq F(b), b \in [0, 1] .$$

□

## C.2 Implement Algorithm 1 Conservatively with Constructed $\bar{F}$ in Exploration Phase

Given the optimistic function  $\bar{F}$ , which is the estimated CDF of the maximum competing bid defined in Eqn. (7). The allocation-payment curve  $\bar{G}$  defined in Definition 3.1 is constructed as  $\bar{G}(b \cdot \bar{F}(b)) = \bar{F}(b)$  for all  $b \in [0, 1]$ .

For each interval  $b \in [c_k, c_{k+1})$ ,  $x$  increases linearly with  $b$  from  $c_k \cdot \bar{F}(c_k)$  to  $c_{k+1} \cdot \bar{F}(c_k)$ , while  $y$  remains constant at  $\bar{F}(c_k)$ . This leads to a "staircase" shape for the  $\bar{G}$  curve when considering all intervals.

The algorithm then computes the concave envelope of this  $\bar{G}$  curve. Due to the attributes of the concave envelope, the concave envelope of  $\bar{G}$  is formed by the linear segments connecting a subset of "breakpoints":  $(c_k \cdot \bar{F}(c_k), \bar{F}(c_k))$  for  $k = 0, 1, \dots, K$ .

Specifically, for a segment between consecutive breakpoints  $(c_k \cdot \bar{F}(c_k), \bar{F}(c_k)), (c_{k+1} \cdot \bar{F}(c_{k+1}), \bar{F}(c_{k+1}))$ , the linear function is given by

$$y_{\text{conv}} = m_i(x_{\text{conv}} - c_k \cdot \bar{F}(c_k)) + \bar{F}(c_k) ,$$

where the slope  $m_k$  will be

$$m_k = \frac{\bar{F}(c_{k+1}) - \bar{F}(c_k)}{c_{k+1} \cdot \bar{F}(c_{k+1}) - c_k \cdot \bar{F}(c_k)} .$$

The algorithm iteratively constructs the upper concave envelope by checking the slopes. If  $m_i < m_{i-1}$ , it indicates a concavity violation at the intermediate point  $(c_k \cdot \bar{F}(c_k), \bar{F}(c_k))$ , and this point is bypassed by forming a new linear segment between  $(c_{k-1} \cdot \bar{F}(c_{k-1}), \bar{F}(c_{k-1}))$  and  $(c_{k+1} \cdot \bar{F}(c_{k+1}))$ .

When  $\bar{F}$  is used as input of Algorithm 1, the algorithm will effectively compute this concave envelope and the corresponding randomized strategy, which will be a probability distribution over

the bids corresponding to the breakpoints of this envelope. These bids are precisely the subset of the grid points  $\{c_k\}_{k=0}^K$  used in the exploration phase.

Given a bid  $b$  generated by Algorithm 1, the conservative strategy, denoted by  $\text{ALG1}^c$ , submits a bid as  $\min(b + 1/K, 1)$ .

**Lemma C.3.** *Given the the optimistic CDF  $\bar{F}$  (defined in Eqn. (7)) and its conservative version  $\bar{F}^c$  (defined in Eqn. (8)). Let  $\text{ALG1}(\bar{F}^c)$  denote the randomized strategy generated by Algorithm 1 with input  $\bar{F}^c$ , and  $\text{ALG1}^c(\bar{F}^c)$  be the corresponding conservative strategy. Fix a value sequence  $V^T$  and run the conservative strategy  $\text{ALG1}^c$  for  $T$  rounds. The total reward of the conservative strategy under the true distribution  $\bar{F}$  satisfies:*

$$\text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) = \text{Reward}(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)).$$

The ROI violation satisfies

$$\Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) \leq \Delta(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)) + \frac{T}{K}.$$

*Proof of Theorem C.3.* Let  $b_t$  be the bid generated by  $\text{ALG}(\bar{F}^c)$  at round  $t$ , and  $b_t^c = \min(b_t + 1/K, 1)$  is the bid of  $\text{ALG1}^c(\bar{F}^c)$ . By the definition of the conservative strategy  $\text{ALG1}^c$ , the output of  $\text{ALG1}^c$

$$b^c = \min(b + 1/K, 1).$$

By the definition of  $\bar{F}^c(b_t)$

$$\bar{F}^c(b_t) = \bar{F}(\min(b_t + 1/K, 1)) = \bar{F}(b_t^c),$$

thus, we can derive that

$$\begin{aligned} \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) &= \sum_{t \in [T]} v_t \cdot \bar{F}(b_t^c) \\ &= \sum_{t \in [T]} v_t \cdot \bar{F}^c(b_t) \\ &= \text{Reward}(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)). \end{aligned}$$

For the payment,

$$\begin{aligned} \text{Payment}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) &= \sum_{t \in [T]} b_t^c \cdot \bar{F}(b_t^c) \\ &= \sum_{t \in [T]} \min(b_t + 1/K, 1) \cdot \bar{F}^c(b_t) \\ &\leq \text{Payment}(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)) + \frac{T}{K}. \end{aligned}$$

Thus, the ROI violation satisfies

$$\begin{aligned} \Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) &= \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) - \text{Payment}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) \\ &\leq \text{Reward}(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)) - \text{Payment}(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)) + \frac{T}{K} \\ &\leq \Delta(\text{ALG1}(\bar{F}^c) | (\bar{F}^c, V^T)) + \frac{T}{K}. \end{aligned} \quad \square$$

### C.3 Analyze the Performance of Implementing Algorithm 1 Conservatively

The optimal randomized strategy algorithm (Algorithm 1), discussed in Section 3 (with a known true distribution  $F$ ), achieves a regret of  $\sqrt{T}$  and a ROI violation of  $O(\sqrt{T} \cdot \log T)$  relative to the optimal randomized strategy subject to the strict ROI constraint under  $F$ .

In Section C.2, we detailed how the randomized bidding strategy is designed using the optimistic distribution  $\bar{F}$  constructed during the exploration phase. The bids submitted are drawn from the set of bids explored.

This section provides a theoretical analysis of the conservative strategy based on the output of Algorithm 1, which was introduced in Section C.2. We aim to understand its performance in terms of reward and ROI violation when the underlying distribution of competing bids is unknown and we rely on our estimate  $\bar{F}^c$ .

**Lemma C.4.** *Consider a fixed value sequence  $V^T$  and a true CDF  $F$  of the maximum competing bid. The optimistic CDF  $\bar{F}$  (Eqn. (7)) and its conservative variant  $\bar{F}^c$  (Eqn. (8)) are obtained in the exploration phase of Algorithm 3 using a grid point set  $\{c_k\}_{k=0}^K (c_k = k/K)$ . For points in the point set, it satisfies  $\bar{F}(c_k) \geq F(c_k)$  and  $\bar{F}(c_k) - F(c_k) \leq 2\varepsilon, \forall c_k \in \{c_k\}_{k=0}^K$ . Let  $\text{ALG1}^c(\bar{F}^c)$  denote the conservative strategy based on Algorithm 1 with input  $\bar{F}^c$ . Then, Running  $\text{ALG1}^c(\bar{F}^c)$  for  $T$  rounds, the total reward over the true distribution  $F$  satisfies:*

$$\text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \geq \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) - T \cdot 2\varepsilon,$$

and the ROI violation cannot be amplified by more than  $T \cdot 2\varepsilon$ .

$$\Delta(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \leq \Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) + T \cdot 2\varepsilon.$$

*Proof of Theorem C.4.* Consider a single round  $t$ , where the input value is  $v_t$ . Algorithm 1  $\text{ALG}(\bar{F}^c)$  generates an optimal bid  $b_{\text{conv}}$ , which, as analyzed in Theorem 3.1, obtains a randomized strategy  $b_1$  with probability  $p_1$  and  $b_2$  with probability  $p_2$ .  $p_1 + p_2 = 1$  and  $p_1 \geq 0, p_2 \geq 0$ . Section C.2 showed that  $b_1$  and  $b_2$  are in the grid points  $\{c_k\}_{k=0}^K (c_k = k/K)$ . Further, the conservative strategy  $\text{ALG1}^c$  will generate  $b_1^c = \min(b_1 + 1/K, 1), b_2^c = \min(b_2 + 1/K, 1)$ , which are still in the point set  $\{c_k\}_{k=0}^K$ . Under  $\bar{F}$ , the expected reward for the conservative strategy in this round, denoted by  $\bar{r}_t$ , satisfies:

$$\bar{r}_t = v_t \cdot [p_1 \cdot \bar{F}(b_1^c) + p_2 \cdot \bar{F}(b_2^c)].$$

The expected payment, denoted by  $\bar{p}_t$ , satisfies:

$$\bar{p}_t = p_1 \cdot b_1^c \cdot \bar{F}(b_1^c) + p_2 \cdot b_2^c \cdot \bar{F}(b_2^c).$$

Similarly, under  $F$ , the expected reward, denoted by  $r_t$ , satisfies:

$$r_t = v_t \cdot [p_1 \cdot F(b_1^c) + p_2 \cdot F(b_2^c)].$$

The expected payment, denoted by  $p_t$ , satisfies:

$$p_t = p_1 \cdot b_1^c \cdot F(b_1^c) + p_2 \cdot b_2^c \cdot F(b_2^c).$$

Given  $|\bar{F}(b) - F(b)| \leq 2\varepsilon$  and  $F(b) \leq \bar{F}(b)$ , for any  $b \in \{c_k\}_{k=0}^K$ , the difference in value would be  $F(b) \geq \bar{F} - 2\varepsilon$ . Difference in expected reward in round  $t$  satisfies:

$$\begin{aligned} r_t - \bar{r}_t &= v_t \cdot \{p_1 \cdot [F(b_1^c) - \bar{F}(b_1^c)] + p_2 \cdot [F(b_2^c) - \bar{F}(b_2^c)]\} \\ &\geq v_t \cdot [p_1 \cdot (-2\varepsilon) + p_2 \cdot (-2\varepsilon)] \\ &\geq -2\varepsilon. \end{aligned}$$

Difference in expected payment in round  $t$  satisfies:

$$p - \bar{p} = p_1 \cdot b_1^c \cdot [F(b_1^c) - \bar{F}(b_1^c)] + p_2 \cdot b_2^c \cdot [F(b_2^c) - \bar{F}(b_2^c)] \leq 0 .$$

Therefore, over  $T$  rounds, the expected reward satisfies:

$$\begin{aligned} \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) &= \sum_{t \in [T]} r_t \\ &\geq \sum_{t \in [T]} \bar{r}_t - T \cdot 2\varepsilon \\ &= \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) - T \cdot 2\varepsilon . \end{aligned}$$

Similarly, the ROI violation over  $T$  rounds satisfies:

$$\begin{aligned} \Delta(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) &= \text{Payment}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) - \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \\ &\leq \text{Payment}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) - \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) + T \cdot 2\varepsilon \\ &= \Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) + T \cdot 2\varepsilon . \end{aligned} \quad \square$$

We analyze the Regret and ROI violations of the conservative strategy based on Algorithm 1 with an estimated distribution against the Optimal strategy under true distribution.

**Lemma C.5.** *Consider a fixed value sequence  $V^T$  and a true CDF  $F$  of the maximum competing bid. Suppose that after the exploration phase of Algorithm 3, our empirical estimate  $\hat{F}$  satisfies  $|F(c_k) - \hat{F}(c_k)| \leq \varepsilon$ , for all grid points  $\{c_k\}_{k=0}^K (c_k = k/K)$ . We define the optimistic distribution  $\bar{F}(b) = \min(\hat{F}(b) + \varepsilon, 1)$  and its conservative version  $\bar{F}^c(c_k) = \bar{F}(\min(c_k + 1/K, 1))$ . Then, run the conservative strategy proposed in Section C.2 with  $\bar{F}^c$  as input, denoted as  $\text{ALG1}^c(\bar{F}^c)$ , for  $T$  rounds. The total reward satisfies:*

$$\text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \geq \text{Reward}(\text{OPT}(F) | (F, V^T)) - O(\sqrt{T}) - T \cdot 2\varepsilon .$$

The ROI violation is bounded by:

$$\Delta(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \leq 2\sqrt{T} \log T + T \cdot 2\varepsilon + \frac{T}{K} .$$

*Proof of Theorem C.5.* Given  $|F(c_k) - \hat{F}(c_k)| \leq \varepsilon$ , and  $\bar{F}(c_k) = \min(\hat{F}(c_k) + \varepsilon, 1)$ , we have  $F(c_k) \leq \bar{F}(c_k)$ . Also,  $|\bar{F}(c_k) - F(c_k)| \leq 2\varepsilon$  for all grid points  $\{c_k\}_{k=0}^K$ .

From Theorem C.4, we have:

$$\text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \geq \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) - T \cdot 2\varepsilon .$$

From Theorem C.3, we know:

$$\text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) = \text{Reward}(\text{ALG}(\bar{F}^c) | (\bar{F}^c, V^T)) .$$

Theorem 3.2 states that the Regret of Algorithm 1 is  $O(\sqrt{T})$ :

$$\text{Reward}(\text{ALG}(\bar{F}^c) | (\bar{F}^c, V^T)) = \text{Reward}(\text{OPT}(\bar{F}^c) | (\bar{F}^c, V^T)) - O(\sqrt{T}) .$$

Since Theorem C.2 showed that  $\bar{F}^c(b) \geq F(b)$  over the entire range  $b \in [0, 1]$  and Theorem B.3 implies that the optimal reward of  $\bar{F}^c$  will not be worse than  $F$ ,

$$\text{Reward}(\text{OPT}(\bar{F}^c) | (\bar{F}^c, V^T)) \geq \text{Reward}(\text{OPT}(F) | (F, V^T)) .$$

Combining these, we have:

$$\begin{aligned} \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) &\geq \text{Reward}(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) - T \cdot 2\varepsilon \\ &= \text{Reward}(\text{ALG}(\bar{F}^c) | (\bar{F}^c, V^T)) - T \cdot 2\varepsilon \\ &= \text{Reward}(\text{OPT}(\bar{F}^c) | (\bar{F}^c, V^T)) - T \cdot 2\varepsilon - O(\sqrt{T}) \\ &\geq \text{Reward}(\text{OPT}(F) | (F, V^T)) - T \cdot 2\varepsilon - O(\sqrt{T}) . \end{aligned}$$

Similarly, for the ROI violation, from Theorem C.4, we have:

$$\Delta(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) \leq \Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) + T \cdot 2\varepsilon .$$

From Theorem C.3, we have:

$$\Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) \leq \Delta(\text{ALG}(\bar{F}^c) | (\bar{F}^c, V^T)) + \frac{T}{K} .$$

From Theorem 3.2, we have:

$$\Delta(\text{ALG}(\bar{F}) | (\bar{F}, V^T)) \leq 2\sqrt{T} \log T .$$

Combining above, the ROI violation is bounded by:

$$\begin{aligned} \Delta(\text{ALG1}^c(\bar{F}^c) | (F, V^T)) &\leq \Delta(\text{ALG1}^c(\bar{F}^c) | (\bar{F}, V^T)) + T \cdot 2\varepsilon \\ &\leq \Delta(\text{ALG}(\bar{F}^c) | (\bar{F}^c, V^T)) + T \cdot 2\varepsilon + \frac{T}{K} \\ &\leq 2\sqrt{T} \log T + T \cdot 2\varepsilon + \frac{T}{K} . \end{aligned} \quad \square$$

#### C.4 Regrets and ROI Violation Bound of the Algorithm 3

This section derives the total regret and ROI violation of the proposed Algorithm 3 by considering both the exploration and exploitation phases.

**Lemma C.6.** *Consider an unknown CDF  $F$  of maximum competing bid. Given a fixed value sequence  $V^T$ , run Algorithm 3 for  $T$  rounds. We sample  $M$  times for each grid point in the point set  $\{c_k\}_{k=0}^{K-1}$  in the exploration phase. The reward of Algorithm 3 satisfies*

$$\text{Regret}(\text{Algorithm 3} | (F, V^T)) \leq M \cdot K + O(\sqrt{T}) + T \cdot 2\varepsilon . \quad (16)$$

The ROI violation satisfies

$$\Delta(\text{Algorithm 3} | (F, V^T)) \leq M \cdot K + 2\sqrt{T} \log T + T \cdot 2\varepsilon + \frac{T}{K} , \quad (17)$$

where  $\varepsilon$  is a parameter such that the empirical CDF  $\hat{F}$  obtained in the exploration phase satisfies:  
 $\Pr[\exists i \in [K] : |\hat{F}(c_k) - F(c_k)| \geq \varepsilon] \leq K \cdot \exp(-2M\varepsilon^2)$ .

*Proof of Theorem C.6.* The algorithm operates in two phases. In the exploration phase, lasting  $T_1 = M \cdot K$  rounds, we bid at  $K$  different prices. In the worst case, the regret in each of these rounds could be as high as the potential reward, which is bounded by 1. Thus, the regret of the exploration phase, denoted by  $\text{Regret}_1$ , is bounded by

$$\text{Regret}_1 \leq \text{Reward}(\text{OPT}(F) | (F, V^{T_1})) \leq \sum_{i \in [M \cdot K]} v_i \leq M \cdot K ,$$

where  $V^{T_1}$  is the value sequence in the exploration phase. Similarly, the total payment, denoted by  $\text{Payment}_1$ , is bounded by

$$\text{Payment}_1 = \sum_{t \in [T_1]} b_t \cdot F(b_t) \leq T_1 \leq M \cdot K .$$

In the worst case for the ROI violation, it will be as high as payment:

$$\Delta_1 \leq \text{Payment}_1 - \text{Reward}_1 \leq T_1 \leq M \cdot K .$$

In the exploitation phase, lasting  $T_2 = 1 - M \cdot K$  rounds, we employ the conservative strategy  $\text{ALG1}^c$  based on the conservative version of the optimistic CDF  $\bar{F}^c$ . From Theorem C.5 we obtain that the reward of the conservative strategy satisfies:

$$\begin{aligned} \text{Reward}(\text{ALG1}^c(\bar{F}) | (F, V^{T_2})) &\geq \text{Reward}(\text{OPT}(F) | (F, V^{T_2})) - o(\sqrt{T_2}) - T_2 \cdot 2\epsilon \\ &\geq \text{Reward}(\text{OPT}(F) | (F, V^{T_2})) - O(\sqrt{T}) - T \cdot 2\epsilon , \end{aligned}$$

where  $V^{T_2}$  is the value sequence in the exploitation phase.

The regret of the exploitation phase, denoted by  $\text{Regret}_2$ , is bounded by

$$\text{Regret}_2 = \text{Reward}(\text{OPT}(F) | (F, V^{T_2})) - \text{Reward}(\text{ALG1}^c(\bar{F}) | (F, V^{T_2})) \leq O(\sqrt{T}) + T \cdot 2\epsilon .$$

The ROI violation of the exploitation phase, denoted by  $\Delta_2$ , satisfies

$$\Delta_2 \leq 2\sqrt{T_2} \log T_2 + T_2 \cdot 2\epsilon + \frac{T_2}{K} .$$

Combining the regret from both phases, we have:

$$\text{Regret}(\text{Algorithm 3} | (F, V^T)) \leq M \cdot K + O(\sqrt{T}) + T \cdot 2\epsilon .$$

Combining the ROI violation from both phases, we have:

$$\Delta(\text{Algorithm 3} | (F, V^T)) \leq M \cdot K + 2\sqrt{T} \log T + T \cdot 2\epsilon + \frac{T}{K} . \quad \square$$

## C.5 Put All Pieces Together to Prove Theorem 5.1

*Proof of Theorem 5.1.* We first discuss the regret and ROI violation incurred with any fixed value sequence  $V^T$ . The proposed algorithm splits into 2 phases: Exploration and Exploitation. In the exploration phase, we identify  $K + 1$  grid points  $c_k = k/K$  for  $k = 0, \dots, K$ . We then sample the points  $\{c_k\}_{k=0}^{K-1}$   $M$  times each, which leads to  $M \cdot K$  rounds. The point  $c_K = 1$  is generally not sampled as its outcome  $F(1) = 1$  is assumed known.

By Theorem C.6, The total Regret combining two phases is bounded by

$$\text{Regret}(\text{Algorithm 3} | (F, V^T)) \leq M \cdot K + O(\sqrt{T}) + T \cdot 2\varepsilon ,$$

with ROI violation:

$$\Delta(\text{Algorithm 3} | (F, V^T)) \leq M \cdot K + 2\sqrt{T} \log T + T \cdot 2\varepsilon + \frac{T}{K} ,$$

where  $\varepsilon$  is the parameter to measure the accuracy of the empirical CDF  $\hat{F}$  we constructed in Section C.1 relative to the true CDF  $F$ . From Theorem C.1, they are bounded by

$$\Pr\left[\exists k \in [K] : |\hat{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq K \cdot \exp(-2M\varepsilon^2) .$$

For each grid point, it satisfies:

$$\Pr\left[|\hat{F}(c_k) - F(c_k)| \geq \varepsilon\right] \leq \exp(-2M\varepsilon^2) ,$$

Denote  $\exp(-2M\varepsilon^2)$  by  $\delta$ . We need to choose  $\varepsilon$  such that the probability of large deviation in the empirical CDF is small. Let's aim for a failure probability  $\delta = T^{-c}$ ,  $c > 0$ . To guarantee that the union bound will decrease,

$$K = o\left(\frac{1}{\delta}\right) = o(T^\alpha) .$$

From Theorem C.6

$$\begin{aligned} \Delta(\text{Algorithm 3} | (F, V^T)) &\leq M \cdot K + 2\sqrt{T} \log T + T \cdot 2\varepsilon + \frac{T}{K} \\ &\leq M \cdot T^\alpha + 2\sqrt{T} \log T + T \cdot 2\sqrt{\alpha \cdot \frac{\log T}{2M}} + T^{1-\alpha} . \end{aligned}$$

Therefore, The ROI violation will be bounded by one of these:  $M \cdot T^\alpha$ ,  $T \cdot 2\sqrt{\alpha \cdot \frac{\log T}{2M}}$ , and  $T^{1-\alpha}$ . The minimum bound will be obtained when

$$O(M \cdot T^\alpha) = O\left(\frac{T}{\sqrt{M}}\right) = O(T^{1-\alpha}) .$$

Thus, with choosing  $M = T^{1/2}$ ,  $K = T^{1/4}$ , the ROI violation can be bounded as:

$$\begin{aligned} \Delta(\text{Algorithm 3} | (F, V^T)) &\leq T^{3/4} + 2\sqrt{T} \log T + T^{3/4} \cdot 2\sqrt{\alpha \cdot \log T} + T^{3/4} \\ &= \tilde{O}(T^{3/4}) . \end{aligned}$$

Substitute  $M = T^{1/2}$ ,  $K = T^{1/4}$  into the regret bound:

$$\begin{aligned} \text{Regret}(\text{Algorithm 3} | (F, V^T)) &\leq M \cdot K + O(\sqrt{T}) + T \cdot 2\varepsilon \\ &\leq T^{3/4} + O(\sqrt{T}) + T^{3/4} \cdot 2\sqrt{\alpha \cdot \log T} \\ &= \tilde{O}(T^{3/4}) . \end{aligned}$$

Thus, when a value sequence  $V^T$  is i.i.d drawn from  $H^T$ , the expected Regret and ROI violation over the distribution  $F$  and  $H$  satisfies

$$\begin{aligned} \text{Regret}(\text{Algorithm 3} | (F, H)) &= \mathbb{E}_{V^T \sim H^T} [\text{Regret}(\text{Algorithm 3} | (F, V^T))] = \tilde{O}(T^{3/4}) , \\ \Delta(\text{Algorithm 3} | (F, H)) &= \mathbb{E}_{V^T \sim H^T} [\Delta(\text{Algorithm 3} | (F, V^T))] = \tilde{O}(T^{3/4}) . \end{aligned} \quad \square$$