

Leveraging Reviews: Learning to Price with Buyer and Seller Uncertainty

Wenshuo Guo*
UC Berkeley
wsguo@berkeley.edu

Nika Haghtalab
UC Berkeley
nika@berkeley.edu

Kirthivasan Kandasamy
UW Madison
kandasamy@cs.wisc.edu

Ellen Vitercik
Stanford University
vitercik@stanford.edu

September 12, 2023

Abstract

In online marketplaces, customers have access to hundreds of reviews for a single product. Buyers often use reviews from other customers that share their type—such as height for clothing, skin type for skincare products, and location for outdoor furniture—to estimate their values, which they may not know *a priori*. Customers with few relevant reviews may hesitate to make a purchase except at a low price, so for the seller, there is a tension between setting high prices and ensuring that there are enough reviews so that buyers can confidently estimate their values. Simultaneously, sellers may use reviews to gauge the demand for items they wish to sell.

In this work, we study this pricing problem in an online learning setting where the seller interacts with a set of buyers of finitely many types, one by one, over a series of T rounds. At each round, the seller first sets a price. Then a buyer arrives and examines the reviews of the previous buyers with the same type, which reveal those buyers' *ex-post* values. Based on the reviews, the buyer decides to purchase if they have good reason to believe that their *ex-ante* utility is positive. Crucially, the seller does not know the buyer's type when setting the price, nor even the distribution over types. We provide a no-regret algorithm that the seller can use to obtain high revenue. When there are d types, after T rounds, our algorithm achieves a problem-independent $\tilde{O}(T^{2/3}d^{1/3})$ regret bound. However, when the smallest probability q_{\min} that any given type appears is large, specifically when $q_{\min} \in \Omega(d^{-2/3}T^{-1/3})$, then the same algorithm achieves a $\tilde{O}(T^{1/2}q_{\min}^{-1/2})$ regret bound. Our algorithm starts by setting lower prices initially so as to (i) boost the number of reviews and increase the accuracy of future buyers' value estimates while also (ii) allowing the seller to identify which customers need to be targeted to maximize revenue. This mimics real-world pricing dynamics. We complement these upper bounds with matching lower bounds in both regimes, showing that our algorithm is minimax optimal up to lower-order terms.

*Work done while the author was a PhD student at UC Berkeley.

1 Introduction

The rapid growth of e-commerce, now accounting for 22% of global retail sales¹, has allowed customers to make far more informed purchase decisions than ever before. Potential buyers can gain insights from thousands of reviews before deciding whether to purchase an item. Customers often use reviews by buyers who share their “type”—such as body type for clothes or skin type for skin-care products—to develop high-fidelity estimates of how much they value different items, which are quantities they may be uncertain of before purchasing.

When learning from reviews, a customer’s purchase decision is no longer just a function of the item’s price but also of how certain the customer is about her valuation, which in turn depends on the earlier sales and reviews of the items. This leads to a tension between setting revenue-optimal prices while ensuring that buyers have enough reviews to confidently estimate their values. This tension is perhaps most clear for customers of rare types (for example, particularly tall or short individuals shopping for clothing) who may find only a few reviews from similar customers and, due to this uncertainty, may only be willing to buy at relatively low prices.

We introduce a model that simultaneously captures the seller’s pricing problem, the buyers’ learning problem, and the modus through which the buyers learn: reviews. We study how a seller—who is uncertain about the buyers’ type distribution—can learn to set high-revenue prices when the buyers themselves are uncertain about their own values and are learning from reviews. Thus, there is information uncertainty on both sides of the market: the seller has uncertainty about which buyer will arrive and the buyers’ type distribution, but the buyer, who knows their type, suffers from the uncertainty about their *ex-ante* value. Both sides of the market are operating with significantly less information than has historically been assumed in mechanism design. We study this pricing problem with an online sequential learning model where the seller attempts to sell identical copies of an item to a series of distinct buyers over T timesteps. Each buyer has one of d types drawn from a distribution \mathcal{P} , and a buyer of type i has an *ex-ante* value of θ_i for the item.

At each timestep t , the seller sets a price p_t . Although the seller knows the *ex-ante* values $\theta_1, \dots, \theta_d$ and thus has some limited information about the buyers (for example, from market research), he does not know the buyer’s type on each round nor even the distribution \mathcal{P} . A buyer on any round could be of (i) a high-value type, but who is uncertain of their value since their type has few reviews, and thus may be hesitant to make a purchase except at a low price, (ii) of a high-value type, and who is more certain of their value since their type has many reviews, and thus is willing to purchase at a high price, or (iii) of a low-value type whom the seller should not target even if they were absolutely certain of their value since it leads to small per-purchase revenue.

If a buyer of type i purchases the item, they will leave a review communicating their *ex-post* value for the item, which is a random variable with mean θ_i . To decide whether to purchase, a new buyer evaluates reviews left by buyers of type i who bought the item in the past. Specifically, the buyer at round $t \in [T]$ uses the past reviews to select a threshold τ_t and chooses to buy as long as $p_t \leq \tau_t$. If the buyer’s threshold τ_t is too pessimistic—for example, it always equals zero no matter the reviews—then optimizing revenue would be hopeless. In our model, we bound the level of pessimism that the buyer can display: we assume that τ_t is at least a lower confidence bound we denote LB_t that equals the average of the reviews left by buyers with the same type, minus an uncertainty term that depends on the number of such reviews. Intuitively, the buyer can be

¹<https://www.trade.gov/ecommerce-sales-size-forecast>

confident that their *ex-ante* value is at least LB_t with high probability, so they always buy if they have good reason to believe that their *ex-ante* utility (value minus price) will be positive.

The *ex-post* value is the actual experience of the buyer and is different from the *ex-ante* value due to exogenous stochastic factors that cannot be known at the time of purchase (for example, manufacturing defects, color on the website not matching the actual color). Hence, the buyer decides based on their *ex-ante* value when there is complete information. In our problem, the buyer does not even know their *ex-ante* value and uses reviews from previous buyers (whose reviews are based on their actual experiences, i.e., *ex-post* values) to update their estimate of the *ex-ante* value (as the expected *ex-post* value is the *ex-ante* value).

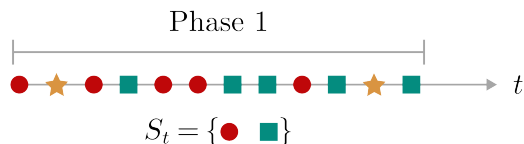
1.1 Our contributions

We provide a no-regret learning algorithm for the seller that balances setting high-revenue prices with soliciting reviews from rare but high-value customers.

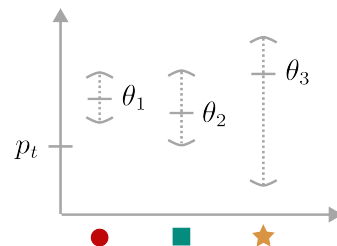
Key technical challenges. The seller does not know the current buyer’s type on each round *a priori*, which means the prices are anonymous. Moreover, this means the seller does not know the number of reviews that the buyer will use to construct their value estimate. If the buyer on round t has a rare type, then the lower confidence bound LB_t will be low, and thus the seller would have to set a low price to ensure a purchase and a review. Suppose this rare type of buyer’s *ex-ante* value is high enough. In that case, it may be worthwhile to initially set a low price to solicit enough reviews to ensure future purchases at a higher price, thereby winning over these rare but high-value customers. The seller, however, has to decide which buyers to win over without knowing the type of the buyer on each round, nor even the distribution over types (and, thus, which types are common and which are rare). He may, therefore, wastefully offer a low price to a high-value buyer with a common type—meaning that LB_t is near the buyer’s *ex-ante* value—who would be willing to buy at a higher price. If a rare buyer’s value is high enough, it may be worthwhile to set a low price to ensure future purchases at a higher price. However, if the buyer’s type is exceedingly rare, the seller will lose too much revenue by setting such a low price. The challenge is that the seller has to decide which buyers to target without knowing the distribution over types.

Algorithm overview. With this intuition in mind, our algorithm maintains a set S_t at each step t consisting of buyer types with a sufficiently high value that are not exceedingly rare. It gradually refines this set over the T rounds. Intuitively, S_t is the set of buyers the algorithm targets. To refine S_t , the algorithm has two phases. In the first phase, the algorithm offers the item for free for a carefully chosen number of rounds, observing i.i.d. samples from the type distribution. The algorithm sets S_t to be the set of types appearing in a sufficiently large fraction of rounds, as in Figure 1a. In the second phase, the algorithm sets the price low enough to ensure that buyers in S_t always buy the item, as in Figure 1b. It successively eliminates types from S_t that contribute too little revenue.

Regret upper bound and proof overview. In this model, we define regret as the difference between (1) the algorithm’s total expected revenue and (2) the expected revenue of the optimal fixed price if the buyers bought whenever their *ex-ante* value was larger than the price, i.e., $\max p \Pr_{i \sim \mathcal{P}}[\theta_i \geq p]$.



(a) Illustration of our algorithm’s first phase, at the end of which only the red circle and the green square are in S_t .



(b) During the second phase, the algorithm sets the price p_t low enough to ensure types in S_t will buy.

Figure 1: Illustration of our algorithm’s first and second phases when there are three types: a red circle, green square, and orange star.

We contend with several sources of regret. The first phase of the algorithm, where the item is sold for free, inevitably leads to regret, so it must be made as brief as possible. The algorithm then completely disregards the buyer types that appeared too rarely during that phase. This results in a subset $Q \subseteq [d]$ of buyer types that appear sufficiently often. In the second phase, the algorithm only attempts to optimize revenue with respect to the buyers in Q instead of the entire set $[d]$, which contributes to regret. Finally, the buyers themselves do not know their *ex-ante* values, whereas, under our regret benchmark, buyers buy whenever their *ex-ante* value is larger than the price.

We obtain our final regret bound by analyzing these three sources of error. Our bound depends on the smallest probability that any given type appears, which we denote as q_{\min} . If q_{\min} is not tiny—specifically, $q_{\min} > 2d^{-2/3}T^{-1/3}$ —then we obtain a regret bound that scales with \sqrt{T} , as desired. In particular, our regret upper bound is $\tilde{O}(T^{1/2}q_{\min}^{-1/2} + T^{1/3}d^{2/3})$. Otherwise, for arbitrary q_{\min} , our regret bound scales with $T^{2/3}$ as $\tilde{O}(T^{2/3}d^{1/3} + T^{1/3}d^{2/3})$.

Regret lower bound and proof overview. Typical bandit lower bounds rely on hypothesis testing arguments to show that any algorithm would struggle to distinguish between similar problems but with different optimal outcomes. Such an analysis would not capture the main difficulty in our setting: how fast customers can estimate their *ex-ante* values from past reviews. Instead, our proof leverages the buyers’ uncertainty to establish a $\tilde{\Omega}(T^{2/3}d^{1/3})$ worst-case lower bound and a $\tilde{\Omega}(T^{1/2}q_{\min}^{-1/2})$ lower bound when q_{\min} is large. This establishes the optimality of our algorithm.

Our proof constructs a hard problem instance where buyer types with low probability of appearance have comparable *ex-ante* values to types with high probability of appearance. On each round, an algorithm should decide whether it wishes to target low-probability customers who may be less certain about their value due to fewer reviews and consequently have small LB_t . Keeping prices low to do so leads to low revenue in the current round, but ignoring low-probability customers by choosing a high price risks losing potentially high per-purchase revenue in the future. By carefully choosing the probability of appearance in our construction, we obtain a tight lower bound.

Our lower bound proof also provides insights that support the structure of our learning algorithm. If the seller knew the type distribution, he could choose a threshold *a priori* and only target customer types with probability larger than that threshold. Our proof illustrates that no policy could do essentially better than this thresholding approach: it does not help significantly to dynamically change which types the seller targets based on appearance probability. Our algorithm

exhibits a similar behavior even though the seller does not know the type distribution: it uses the first phase to discard low-probability types, focusing on the remaining types in the second phase.

1.2 Related work

Learning to price when buyers do not know their values. Learning to price when buyers do not know their values requires new machinery beyond classic pricing algorithms and auction design. Prior works propose different strategies for the seller when the buyers learn through various means. One line of work studies bidding strategies for buyers who do not know their values in auction settings [Feng et al., 2018, Weed et al., 2016, Kandasamy et al., 2023]. Another line of work considers selling repeatedly to a single buyer while the buyer is learning from their own experience at each round [Papadimitriou et al., 2022, Ashlagi et al., 2016, Chawla et al., 2022]. However, a significant limitation in practice is that buyers on online platforms do not necessarily return repeatedly to buy the same item and can only obtain feedback from previous buyers via reviews. In this paper, we study the seller’s pricing strategy when the buyers can only learn from past reviews.

Ifrach et al. [2019] consider a similar pricing problem for the seller when the buyers learn from reviews. However, their model is limited to one buyer type, where the buyers’ values for the item are i.i.d. random variables from a fixed distribution. In contrast, we study the setting where there are multiple buyer types. Moreover, the seller does not know the frequency of each type and the type of buyer who arrives at each round, which leads to crucial difficulties in our analysis.

Learning to price when buyers know their values. Zhao and Chen [2020] study a setting where the buyers know their values, but the seller does not know the distribution over buyers’ values. Reviews give the seller more information about this distribution than purchase decisions alone would. Zhao and Chen [2020] present an algorithm that uses the (non-noisy) reviews to obtain a $\tilde{O}(T^{1/2})$ regret bound. In contrast, if the seller only observes purchase decisions and not reviews, Kleinberg and Leighton [2003] provide a $\Omega(T^{2/3})$ lower bound. While they show that this bound can be improved to $\tilde{\Theta}(T^{1/2})$, it requires additional distributional assumptions.

Selling to no-regret buyers who know their values. In situations where buyers know their values, the buyer may strategically improve their purchase decisions or bidding strategy over repeated interactions to achieve a higher accumulated utility. No-regret learning has been explored as a model of buyer behavior [Braverman et al., 2018, Deng et al., 2019, Nekipelov et al., 2015, Devanur et al., 2014]. In this literature, buyers know their values but may use no-regret algorithms to learn how to bid. In comparison, in this paper, we work with buyers who do not know their values and need to estimate them from historical reviews. This leads to different dynamics. For example, suppose a seller repeatedly sets the Myerson reserve price. In that case, any buyer who knows her value *a priori* and uses a no-regret algorithm will eventually learn to submit a winning bid. However, a buyer without a reasonable estimate of her value may consider the Myerson price too high and will not buy. Interestingly, a seller dealing with either type of learner may benefit from selling the item for a low price early on, but for two very different reasons. In our setting, this will give buyers of a given type the opportunity to refine their estimated value and will encourage future buyers of the same type to buy at higher prices if their value is indeed high. On the other hand, as Braverman et al. [2018] show, giving items for free to agents who are learning to bid will

accrue welfare (as long as agents are allowed to overbid), which the algorithm can then extract in future rounds by setting prices that are higher than agent values.

Buyers’ social learning from reviews. Our work is also related to a rich literature on buyer behavior and social learning from reviews when buyers do not know their values [Ifrach et al., 2019, Boursier et al., 2022, Han and Anderson, 2020, Chamley, 2004, Besbes and Scarsini, 2018, Bose et al., 2006, Crapis et al., 2017, Kakhbod et al., 2021, Acemoglu et al., 2022]. Much of the research on social learning from reviews can be categorized into two groups depending on whether the decision model is Bayesian or non-Bayesian. In the Bayesian model, Ifrach et al. [2019], Acemoglu et al. [2022], and Boursier et al. [2022] study a setting where the buyers decide whether to purchase the item by calculating posterior probabilities about the item’s quality given the past reviews.

It may be computationally challenging for buyers to compute Bayesian updates, so several papers relax this assumption [Crapis et al., 2017, Besbes and Scarsini, 2018]. Besbes and Scarsini [2018], for example, study both fully rational Bayesian buyers and buyers with limited rationality who can only observe the average of the past reviews. Under these two extremes, they analyze the conditions under which buyers can recover a product’s true quality based on their observed feedback. Unlike our paper, the buyers have private signals about the item for sale, influencing their purchase decisions. Our model can be seen as situated between these two extremes because the purchase decisions depend on the average of the past reviews and the number of those reviews. Moreover, whereas Besbes and Scarsini [2018] analyze risk-neutral buyers, we study a form of risk aversion where buyers may not purchase even if the price is below the average reviews.

Unlike this prior research, we do not assume all buyers share a specific decision policy. Instead, we identify a broad family of decision policies under which our results hold. In particular, we only require that the buyer purchases the item if the price is sufficiently low.

2 Notation and online learning setup

In our model, an item is sold repeatedly to a sequence of distinct buyers over a series of T rounds. Each buyer has a type $i \in [d]$, and there is an unknown distribution \mathcal{P} over the types $[d]$. We use the notation $q_i = \Pr_{j \sim \mathcal{P}}[j = i]$ and $q_{\min} = \min_{i \in [d]} q_i$.

The *ex-ante* value of a buyer with type $i \in [d]$ is $\theta_i \in [0, 1]$. If a buyer with type $i \in [d]$ purchases the item, their *ex-post* value is drawn from a distribution \mathcal{D}_i with support $[0, 1]$ and mean θ_i . The seller knows $\theta_1, \dots, \theta_d$ but not the distributions $\mathcal{P}, \mathcal{D}_1, \dots, \mathcal{D}_d$. For ease of analysis, we assume that the seller has ordered the types such that $\theta_1 \leq \theta_2 \leq \dots \leq \theta_d$, but the buyers are unaware of this ordering. This assumption is not necessary for the results to hold.

At each timestep $t \in [T]$:

1. There is a set σ_{t-1} of reviews which describe past buyers’ types and their *ex-post* values.
2. The seller first sets a price $p_t \in [0, 1]$.
3. A buyer arrives with type $i_t \sim \mathcal{P}$. They observe the past reviews of buyers with type i_t : $\Phi_{i_t, t} = \{v : (i, v) \in \sigma_{t-1} \text{ and } i = i_t\}$. They decide whether to purchase the item using $\Phi_{i_t, t}$. We describe the buyer’s purchasing model in more detail in Section 2.1. Observe that the seller is unaware of the buyer’s type i_t when they set the price.

4. If the buyer purchases the item, they pay p_t and leave a review of (i_t, v_t) describing both their type and their *ex-post* value $v_t \sim \mathcal{D}_{i_t}$. In this case, $\sigma_t = \sigma_{t-1} \cup \{(i_t, v_t)\}$, and otherwise, $\sigma_t = \sigma_{t-1}$.

Our assumptions and model reflect practical e-commerce settings. First, quite often, it is reasonable to assume that sellers know customers’ *ex-ante* values as they may have inside information. For instance, a skincare product vendor may know that a particular product works better on some skin types. However, buyers may not simply trust the seller if they were to publish this value, as the seller has every incentive to overstate this value to maximize revenue. A buyer would instead decide if a product is suitable for her via independent reviews from other customers. Second, for fairness reasons, in e-commerce platforms, sellers typically have to publish a single price for all customers and cannot sell the item at individualized prices. Third, if a buyer does not purchase an item, they will not leave a review, and the seller has no way of knowing their type or *ex-post* value.

2.1 Buyers’ purchasing model

At time step t , the agent’s purchase decision is defined by a threshold $\tau_t(\sigma_{t-1}, i_t) \geq 0$ that takes as input their type i_t and the reviews left by past agents. Intuitively, $\tau_t(\sigma_{t-1}, i_t)$ represents the agent’s estimate of their value θ_{i_t} based on past reviews. The agent purchases the item if $p_t \leq \tau_t(\sigma_{t-1}, i_t)$.

A conservative agent would choose $\tau_t(\sigma_{t-1}, i_t)$ to be low in order to always guarantee that $\tau_t(\sigma_{t-1}, i_t) \leq \theta_{i_t}$, so that they only purchase when their *ex-ante* utility is non-negative. An extreme example of this type of conservatism would set $\tau_t(\sigma_{t-1}, i_t) = 0$, meaning that the agent would only purchase the item if offered for free. Optimizing revenue with such a conservative agent would be hopeless. Therefore, we impose the following natural lower bound on $\tau_t(\sigma_{t-1}, i_t)$:

Definition 2.1. Let $\Phi_t \subseteq \sigma_{t-1}$ be the reviews left by agents with type i_t :

$$\Phi_t = \{v : (i, v) \in \sigma_{t-1} \text{ and } i = i_t\}.$$

Let LB_t be the average of these reviews minus a standard confidence term:

$$\text{LB}_t = \begin{cases} 0 & \text{if } \Phi_t = \emptyset, \\ \max \left\{ 0, \frac{1}{|\Phi_t|} \sum_{v \in \Phi_t} v - \sqrt{\frac{1}{2|\Phi_t|} \ln \frac{t}{\eta}} \right\} & \text{else.} \end{cases}$$

We say that the agent on round t is η -*pessimistic* if, $\tau_t(\sigma_{t-1}, i_t) \geq \text{LB}_t$.

This uncertainty term corresponds to the standard confidence interval defined by the Hoeffding bound. Intuitively, as a buyer sees more reviews from his type, this uncertainty decreases, and he is more certain about his *ex-ante* valuation. The $\ln t$ term is necessary to construct a valid confidence interval for an arbitrary algorithm as the data may not be independent (see Appendix A): the algorithm’s price may depend on previous reviews, which in turn will affect future buyers and reviews. This $\ln t$ term is not fundamental—the lower bound does not use it.

Intuitively, the agents can be confident that *regardless* of the policy used by the seller, with probability $1 - \eta$, for all rounds $t \in [T]$, $\theta_{i_t} \geq \text{LB}_t$. We prove this formally in Appendix A. Therefore, if the price is lower than LB_t , an η -pessimistic agent will buy the item as they can be confident, based on past reviews, that their *ex-ante* utility $\theta_{i_t} - p_t$ will be non-negative. This restriction bounds the level of pessimism that the agents can display and thus makes it possible to set reasonable prices. We clip this lower confidence bound at 0 since valuations are always in $[0, 1]$.

Algorithm 1: TYPEELIMINATION

Input: Number of timesteps t_λ , number of types d , parameter $\lambda \in [0, 1]$

for $t = 1, \dots, t_\lambda$ **do**

 Set $p_t = 0$

 Buyer with type i_t arrives and purchases item

 Buyer leaves review (i_t, v_t) , where $v_t \sim \mathcal{D}_{i_t}$

end

for $i \in [d]$ **do**

 Set

$$\bar{q}_i = \frac{1}{t_\lambda} \sum_{t=1}^{t_\lambda} \mathbb{I}(i_t = i)$$

 ▷ Calculate the fraction of rounds each type appeared

end

Set $Q = \{i : \bar{q}_i \geq \frac{3\lambda}{4}\}$ ▷ Set of types that appeared at least a $(3\lambda/4)$ -fraction of rounds

Output: Q

2.2 Regret

We define regret as the difference between:

1. The algorithm's total expected revenue, and
2. (*baseline*) The expected revenue of the optimal fixed price if the agents bought whenever their *ex-ante* value was larger than the price.

Under the baseline that we compete with, both the buyer and the seller are equipped with more information than in the learning problem: the seller knows all distributions $\mathcal{P}, \mathcal{D}_1, \dots, \mathcal{D}_d$ and the buyers know their *ex-ante* values $\theta_1, \dots, \theta_d$. Therefore, the seller knows *a priori* which customers to target to maximize revenue. Moreover, since the buyers do not need to learn their *ex-ante* values from reviews, the seller can extract higher revenue than they could from uncertain buyers who may only buy when the price is likely lower than their *ex-ante* value.

Formally, let $b_t \in \{0, 1\}$ indicate whether or not the buyer bought on round $t \in [T]$ and let $p^* = \operatorname{argmax}_{p \in [0, 1]} p \Pr_{i \sim \mathcal{P}}[\theta_i \geq p]$ be the price with highest expected revenue if the agents bought whenever their *ex-ante* value was larger than the price. Regret is defined as

$$\mathbb{E}[R_T] = T p^* \Pr_{i \sim \mathcal{P}}[\theta_i \geq p^*] - \mathbb{E} \left[\sum_{t=1}^T p_t b_t \right]. \quad (1)$$

3 Online Pricing Algorithm

This section describes our algorithm, which has two phases: Algorithm 1 and 2. It is defined by a parameter $\lambda > 0$. (We will choose $\lambda = d^{-2/3} T^{-1/3}$ to obtain optimal trade-offs).

Our algorithm has two phases. In the first phase (Algorithm 1), the algorithm sets a price of 0 for $t_\lambda = \Theta(\ln(dT)/\lambda)$ rounds. The agent will buy the item at each round since the price is 0 and leave a review. This allows the algorithm to obtain i.i.d. samples from the type distribution

Algorithm 2: Online pricing with reviews

Input: Number of timesteps T , number of types d , $\eta \in [0, 1]$ such that the agents are η -pessimistic, parameter $\lambda \in [0, 1]$

Set $t_\lambda \stackrel{\text{def}}{=} \frac{32 \ln(dT^2)}{\lambda} + 1$

Set $S_{t_\lambda+1} = \text{TYPEELIMINATION}(t_\lambda, d, \lambda)$ ▷ S_t is the set of “active types”

for $t = t_\lambda + 1, \dots, T$ **do**

for $i \in S_t$ **do**

Compute $\Phi_{it} = \{v_s : (i, v_s) \in \sigma_{t-1}\}$ and

$$\text{LB}_{it} = \begin{cases} 0 & \text{if } \Phi_{it} = \emptyset \\ \max \left\{ \frac{1}{|\Phi_{it}|} \sum_{v \in \Phi_{it}} v - \sqrt{\frac{1}{2|\Phi_{it}|} \ln \frac{T}{\eta}}, 0 \right\} & \text{else} \end{cases}$$

end

Set price $p_t = \min_{i \in S_t} \{\min\{\theta_i, \text{LB}_{it}\}\}$ ▷ p_t is the smallest LB_{it} or θ_i of any active type

$b_t = \mathbb{I}(\text{buyer buys at price } p_t)$ ▷ We prove that if $i_t \in S_t$, then $b_t = 1$

if $b_t = 1$ **then**

| The buyer leaves a review (i_t, v_t) where $v_t \sim \mathcal{D}_{i_t}$

end

Set $\rho_t = \sqrt{\frac{\ln(dT^2)}{2(t-t_\lambda)}}$

for $i \in S_t$ **do**

$\bar{\mu}_{i,t} = \frac{1}{t-t_\lambda} \sum_{s=t_\lambda+1}^t \theta_i \cdot \mathbb{I}(b_s = 1 \wedge \theta_{i_s} \geq \theta_i \wedge i_s \in Q)$ ▷ Estimate of rev (θ_i, Q)

$\hat{\mu}_{i,t} = \bar{\mu}_{i,t} + \rho_t$ ▷ Upper confidence bound

$\check{\mu}_{i,t} = \bar{\mu}_{i,t} - \rho_t$ ▷ Lower confidence bound

end

Set $i_0 = \min\{i \in S_t : \hat{\mu}_{i,t} \geq \max_{k \in S_t} \check{\mu}_{k,t}\}$ ▷ For $i < i_0$, rev (θ_i, Q) is likely too small

Set $S_{t+1} = S_t \cap \{i_0, i_0 + 1, \dots, d\}$ ▷ Eliminate types $i < i_0$

end

\mathcal{P} . In phase 2 (Algorithm 2), i.e., the remaining $T - t_\lambda$ rounds, the algorithm will ignore types that appeared too rarely during phase 1—in particular, on fewer than a $(3\lambda/4)$ -fraction of rounds. Intuitively, customers of these types have a low probability of appearance and thus will have more uncertainty about their values due to fewer reviews. The uncertainty term will cause the lower confidence bound LB_t in Definition 2.1 to be small. As the seller will have to choose a low price to target these customers (even if their ex-ante value is large), they may have to forego higher revenue from more frequent customer types. Therefore, it is not worthwhile for the algorithm to target these customers. We use Q to denote the buyer types that appeared on at least a $(3\lambda/4)$ -fraction of rounds.

To describe the algorithm’s second phase, we will use the notation

$$\text{rev}(p, Q) = p \Pr_{i \sim \mathcal{P}} [\theta_i \geq p \text{ and } i \in Q],$$

to denote the expected revenue of a price p restricted to buyers in Q and $p^*(Q) = \arg\max \text{rev}(p, Q)$. In this phase, Algorithm 2 will ignore the extremely rare buyers not in Q and aim to set prices that

compete with $p^*(Q)$. In the analysis, we will show that by competing with $p^*(Q)$, Algorithm 2 also competes with the optimal price p^* .

Observe that $p^*(Q) = \theta_{i_Q}$ for some $i_Q \in Q$. On each round $t > t_\lambda$ of the second phase, Algorithm 2 maintains a set S_t of “active types” such that i_Q is likely in S_t . Algorithm 2 sets the price p_t low enough to ensure that if the current type i_t is in S_t , then the buyer will buy. In particular, we define LB_{i_t} as the largest price the seller can set to ensure a purchase from a buyer of type i . We then set the price p_t to be the smallest $\text{LB}_{i,t}$ or θ_i of any active type $i \in S_t$ (we include θ_i for ease of analysis). If the buyer purchases the item, they leave a review (i_t, v_t) where $v_t \sim \mathcal{D}_{i_t}$.

Next, for each active type $i \in S_t$, the seller estimates $\text{rev}(\theta_i, Q)$. We denote this estimate as $\bar{\mu}_{i,t}$ along with upper and lower confidence bounds $\hat{\mu}_{i,t}$ and $\check{\mu}_{i,t}$. We will describe this estimate more in Section 4. When estimating the revenue for different prices via the averages $\bar{\mu}_{i,t}$, we only use samples from the second phase. Doing so leads to a cleaner analysis, allowing us to separate the randomness in eliminating low probability types to determine the set Q from the randomness of estimating $\text{rev}(\theta_i, Q)$. However, when constructing the lower confidence bound $\text{LB}_{i,t}$ for customers of type i , we use reviews from all rounds. This is to be expected, as customers will use all past reviews when making a purchasing decision.

Algorithm 2 defines

$$i_0 = \min \left\{ i \in S_t : \hat{\mu}_{i,t} \geq \max_{k \in S_t} \check{\mu}_{k,t} \right\}$$

to be the smallest active type such that θ_{i_0} may plausibly be $p^*(Q)$. For all $i < i_0$, the upper confidence bound on $\text{rev}(\theta_i, Q)$ is small ($\hat{\mu}_{i,t} < \max_{k \in S_t} \check{\mu}_{k,t}$), so it is unlikely that $\theta_i = p^*(Q)$. Algorithm 2 concludes round t by eliminating all types $i < i_0$ from the active set.

4 Regret upper bounds

We now state our main upper bounds on regret (Equation (1)).

Theorem 4.1. *Suppose the agents are η -pessimistic. If $q_{\min} \leq 2\lambda$ then*

$$\mathbb{E}[R_T] = O \left(\frac{\ln(dT)}{\lambda} + Td\lambda + \sqrt{T \ln(dT)} + \sqrt{\frac{T}{\lambda} \ln \frac{dT}{\eta}} \right),$$

and if $q_{\min} > 2\lambda$, then

$$\mathbb{E}[R_T] = O \left(\frac{\ln(dT)}{\lambda} + \sqrt{T \ln(dT)} + \sqrt{\frac{T}{q_{\min}} \ln \frac{dT}{\eta}} \right).$$

Theorem 4.1 implies the following corollary for the specific choice of $\lambda = d^{-2/3}T^{-1/3}$.

Corollary 4.2. *Suppose the agents are η -pessimistic. Setting $\lambda = d^{-2/3}T^{-1/3}$, we have that if $q_{\min} \leq 2d^{-2/3}T^{-1/3}$ then*

$$\mathbb{E}[R_T] = O \left(T^{2/3}d^{1/3} + T^{1/3}d^{1/3} \sqrt{\ln \frac{dT}{\eta}} + \sqrt{T \ln(dT)} + T^{1/3}d^{2/3} \ln(dT) \right),$$

and if $q_{\min} > 2d^{-2/3}T^{-1/3}$, then

$$\mathbb{E}[R_T] = O\left(\sqrt{\frac{T}{q_{\min}} \ln \frac{dT}{\eta}} + T^{1/3}d^{2/3} \ln(dT)\right).$$

We note that while the worst-case regret scales with $T^{2/3}$, it improves to \sqrt{T} when all types appear with large enough probability since customers of all types will be able to form accurate estimates of their values quickly. We emphasize that our algorithm and analysis are markedly different from explore-then-commit (ETC) style algorithms in stochastic bandit settings, which share a similar two-phase strategy and have $T^{2/3}$ regret. First, the first ‘explore’ phase of ETC algorithms is much longer (typically $\tilde{O}(T^{2/3})$ rounds) than our Phase 1, which lasts only $\tilde{O}(T^{1/3})$ rounds. ETC algorithms also focus on learning all unknowns in their first phase, while here, its only purpose is to eliminate low probability types. Second, in the ‘commit’ phase of ETC algorithms, typically, no learning is required, while in our second phase, the algorithm is still learning the optimal price. Third, unlike our algorithm, ETC algorithms cannot obtain \sqrt{T} regret even under favorable conditions [Garivier et al., 2016]. Fourth, we reiterate that the $T^{2/3}$ worst-case regret is due to the uncertainty on the buyers’ side, which is a challenge specific to our setting.

We will first provide an overview of our proof, with the full proof to follow in Section 4.1.

Proof sketch of Theorem 4.1. The terms of our regret bounds in Theorem 4.1 arise from the following steps of our analysis. The first phase immediately contributes $O\left(\frac{\ln(dT)}{\lambda}\right)$ to the regret since the item is sold for free during that phase. At the end of the first phase, Algorithm 1 discards the types that appeared too infrequently, resulting in a set Q , and only aims to maximize revenue over Q . When $q_{\min} \leq 2\lambda$, we prove that competing with $p^*(Q)$ rather than $p^*([d])$ contributes $Td\lambda$ to the regret. Meanwhile, when $q_{\min} > 2\lambda$, we prove that with high probability, $Q = [d]$, and thus $p^*(Q) = p^*([d])$, so there is no impact on regret.

In order to gradually learn a price that competes with $p^*(Q)$, Algorithm 2 maintains estimates $\bar{\mu}_{i,t}$ of $\text{rev}(\theta_i, Q) = \theta_i \Pr_{j \sim \mathcal{P}}[\theta_j \geq \theta_i \text{ and } j \in Q]$ for the active types $i \in S_t$. The error of these estimates contributes a factor of $O\left(\sqrt{T \ln(dT)}\right)$ to the regret. This step of the analysis takes some care because we cannot observe at each round t whether or not $i_t \in Q$, provided the buyer did not buy the item. If we were able to observe whether $i_t \in Q$, we could simply set

$$\bar{\mu}_{i,t} = \frac{1}{t - t_\lambda} \sum_{s=t_\lambda+1}^t \theta_i \cdot \mathbb{I}(\theta_{i_s} \geq \theta_i \text{ and } i_s \in Q),$$

and the concentration would follow from a Hoeffding bound. Instead, we set

$$\bar{\mu}_{i,t} = \frac{1}{t - t_\lambda} \sum_{s=t_\lambda+1}^t \theta_i \cdot \mathbb{I}(b_s = 1, \theta_{i_s} \geq \theta_i, \text{ and } i_s \in Q),$$

but nonetheless prove that it is a good estimate of $\text{rev}(\theta_i, Q)$. To do so, we show that for all active types $i \in S_t$ and all rounds $s \leq t$ of Algorithm 2,

$$\mathbb{I}(\theta_{i_s} \geq \theta_i \text{ and } i_s \in Q) = \mathbb{I}(b_s = 1, \theta_{i_s} \geq \theta_i, \text{ and } i_s \in Q), \quad (2)$$

so we can still apply a Hoeffding bound (taking into account that the set S_t is a random variable). If $b_s = 1$, then clearly Equation (2) holds. Otherwise, $i_s \notin S_s$ because any buyer in S_s will always

buy. We show that this means that either $i_s \notin Q$ or—based on the way that types are eliminated from the active sets— $\theta_{i_s} < \theta_{i_t}$, so Equation (2) holds in this case as well.

Finally, the agents themselves are learning as the algorithm progresses, which increases the regret since our benchmark is the expected revenue of the optimal price if the agents buy whenever their *ex-ante* value is larger than the price. When $q_{\min} \leq 2\lambda$, the fact that the agents are learning contributes $O\left(\sqrt{\frac{T}{\lambda}} \ln \frac{dT}{\eta}\right)$ to the regret and when $q_{\min} > 2\lambda$, it contributes $O\left(\sqrt{\frac{T}{q_{\min}}} \ln \frac{dT}{\eta}\right)$. \square

4.1 Proof of the regret upper bound (Theorem 4.1)

In this section, we prove Theorem 4.1. The proof relies on a handful of helper lemmas which we prove in Appendix B.

Proof of Theorem 4.1. In this proof, on each round $t > t_\lambda$, we use the notation $p'_t = \min_{i \in S_t} \theta_i$. We split the regret into five terms as follows:

$$R_T = \sum_{t=1}^T p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d])) - \sum_{t=1}^T p_t b_t = Z_1 + Z_2 + Z_3 + Z_4 + Z_5.$$

The first term $Z_1 = p^*([d])t_\lambda \leq \frac{32 \ln(dT^2)}{\lambda} + 1$ measures the revenue lost from offering the item for free for the first t_λ rounds. The second term

$$\begin{aligned} Z_2 &= \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d])) - \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) \\ &= \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \notin Q) \end{aligned}$$

relates to the revenue lost due to the fact that we only aim to compete with the optimal price for relatively-common types—namely those in Q —as does the third term

$$Z_3 = \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q).$$

The fourth term

$$Z_4 = \sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q)$$

relates the cumulative revenue of the optimal price over Q —that is, $p^*(Q)$ —to the cumulative revenue of the “proxy” price $p'_t = \min_{i \in S_t} \theta_i$. Finally, the last term

$$Z_5 = \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p_t b_t$$

relates the cumulative revenue of the proxy price p'_t to the algorithm’s cumulative revenue. In the following claims, we bound Z_2, Z_3, Z_4 , and Z_5 . The full proofs are in Appendix B.

Claim 4.3. *If $q_{\min} \leq 2\lambda$ then $\mathbb{E}[Z_2] \leq Td\lambda + 1$ and if $q_{\min} > 2\lambda$, then $\mathbb{E}[Z_2] \leq 1$.*

Proof sketch of Claim 4.3. First, we bound Z_2 as follows:

$$Z_2 = \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \notin Q) \leq \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(i_t \notin Q) \leq \sum_{t>t_\lambda} \mathbb{I}(i_t \notin Q).$$

Recall from Algorithm 1 that \bar{q}_i is the fraction of times that type i appears in phase 1 and let \mathcal{G} be the event that for all $i \in [d]$ such that $q_i \geq \lambda$, we have that $\bar{q}_i \geq \frac{3\lambda}{4}$, which means that $i \in Q$. In other words, when \mathcal{G} happens, $[d] \setminus Q \subseteq \{q_i : q_i < \lambda\}$. In Lemma B.1, we prove that $\Pr[\mathcal{G}^c] \leq \frac{1}{T}$, so $\mathbb{E}[Z_2] \leq \mathbb{E}[Z_2 | \mathcal{G}] + T \Pr[\mathcal{G}^c] \leq \mathbb{E}[Z_2 | \mathcal{G}] + 1$.

Next, since Q is a random variable, we condition on it as well:

$$\mathbb{E}[Z_2 | \mathcal{G}] = \sum_{Q' \subseteq [d]} \mathbb{E}[Z_2 | Q = Q', \mathcal{G}] \Pr[Q = Q' | \mathcal{G}].$$

If $[d] \setminus Q' \not\subseteq \{q_i : q_i < \lambda\}$, then $\Pr[Q = Q' | \mathcal{G}] = 0$. For any Q' such that $[d] \setminus Q' \subseteq \{q_i : q_i < \lambda\}$, we prove

$$\mathbb{E}[Z_2 | Q = Q', \mathcal{G}] = \sum_{t>t_\lambda} \sum_{i \notin Q'} \Pr[i_t = i].$$

If $q_{\min} \leq 2\lambda$, then

$$\sum_{t>t_\lambda} \sum_{i \notin Q'} \Pr[i_t = i] \leq \sum_{t>t_\lambda} \sum_{i \notin Q'} \lambda \leq Td\lambda,$$

which implies that $\mathbb{E}[Z_2] \leq Td\lambda + 1$. The case where $q_{\min} > 2\lambda$ follows similarly. \square

Claim 4.4. $\mathbb{E}[Z_3] \leq 0$.

Proof sketch of Claim 4.4. In this proof we bound

$$Z_3 = \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q). \quad (3)$$

We begin by conditioning the first term of Equation (3) on Q since it is a random variable:

$$\begin{aligned} & \mathbb{E} \left[\sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) \right] \\ &= \sum_{Q' \subseteq [d]} \sum_{t>t_\lambda} p^*([d]) \Pr[\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q' | Q = Q'] \Pr[Q = Q'] \\ &\leq \sum_{Q' \subseteq [d]} (T - t_\lambda) p^*(Q') \Pr_{i \sim \mathcal{P}}[\theta_i \geq p^*(Q') \text{ and } i \in Q'] \Pr[Q = Q'], \end{aligned} \quad (4)$$

where the final inequality follows from the definition of $p^*(Q')$.

Next, for the second term of Equation (3),

$$\begin{aligned} & \mathbb{E} \left[\sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q) \right] \\ &= \sum_{Q' \subseteq [d]} \mathbb{E} \left[\sum_{t>t_\lambda} p^*(Q') \mathbb{I}(\theta_{i_t} \geq p^*(Q') \text{ and } i_t \in Q') | Q = Q' \right] \Pr[Q = Q'] \\ &= \sum_{Q' \subseteq [d]} (T - t_\lambda) p^*(Q') \Pr_{i \sim \mathcal{P}}[\theta_i \geq p^*(Q') \text{ and } i \in Q'] \Pr[Q = Q']. \end{aligned}$$

Combined with Equation (4), we have that $\mathbb{E}[Z_3] \leq 0$. \square

Claim 4.5. $\mathbb{E}[Z_4] \leq 5 + 4\sqrt{2T \ln(dT^2)}$.

Proof sketch of Claim 4.5. In this claim, we bound

$$Z_4 = \sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q) \quad (5)$$

where $p'_t = \min_{i \in S_t} \theta_i$. Beginning with the first term of this equation, we prove that

$$\sum_{t>t_\lambda} \mathbb{E}[p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q)] = \sum_{t>t_\lambda} \mathbb{E}[\text{rev}(p^*(Q), Q)]. \quad (6)$$

Moving on to the second term of Equation (5), we prove that for any $t > t_\lambda$,

$$\mathbb{E}[p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q)] = \mathbb{E}[\text{rev}(p'_t, Q)]. \quad (7)$$

Combining Equations (6) and (7), we have that

$$\mathbb{E}[Z_4] \leq \sum_{t>t_\lambda} \mathbb{E}[\text{rev}(p^*(Q), Q) - \text{rev}(p'_t, Q)]. \quad (8)$$

Next, for all $t > t_\lambda$, let \mathcal{B}_t be the event that:

1. $i_Q \in S_t$ and
2. $\text{rev}(p^*(Q), Q) - \text{rev}(p'_t, Q) \leq 4\rho_{t-1}$ (where $\rho_{t_\lambda} = 1$).

Also, let $\mathcal{C}_t = \bigcap_{s=t_\lambda+1}^t \mathcal{B}_s$. In Lemma B.2, we prove that $\Pr[\mathcal{C}_t^c] \leq \frac{1}{T}$. By Equation (8),

$$\mathbb{E}[Z_4] \leq \mathbb{E}\left[\sum_{t>t_\lambda} \text{rev}(p^*(Q), Q) - \text{rev}(p'_t, Q) \mid \mathcal{C}_T\right] + T \Pr[\mathcal{C}_T^c] \leq \sum_{t>t_\lambda} 4\rho_{t-1} + 1$$

which implies the result. \square

Claim 4.6. If $q_{\min} \leq 2\lambda$, then $\mathbb{E}[Z_5] \leq 4\sqrt{\frac{2T}{\lambda} \ln \frac{dT^2}{\eta}} + 3$ and if $q_{\min} > 2\lambda$, $\mathbb{E}[Z_5] \leq 4\sqrt{\frac{T}{q_{\min}} \ln \frac{dT^2}{\eta}} + 2$.

Proof sketch of Claim 4.6. On each round $t > t_\lambda$, recall that

$$\text{LB}_{it} = \begin{cases} 0 & \text{if } \Phi_{it} = \emptyset \\ \max\left\{\frac{1}{|\Phi_{it}|} \sum_{v \in \Phi_{it}} v - \sqrt{\frac{1}{2|\Phi_{it}|} \ln \frac{T}{\eta}}, 0\right\} & \text{else.} \end{cases}$$

Let $j_t = \text{argmin}_{j \in S_t} \theta_j$, so $p'_t = \theta_{j_t}$. We prove that

$$Z_5 \leq \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \wedge i_t \in Q \wedge b_t = 0) + \sum_{t>t_\lambda} (p'_t - p_t) b_t.$$

Since $p_t \leq p'_t$, we have that

$$Z_5 \leq \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \wedge i_t \in Q \wedge b_t = 0) + \sum_{t>t_\lambda} (p'_t - p_t).$$

By definition of the pricing rule, if $i_t \in S_t$, then $b_t = 1$. Therefore, if $b_t = 0$, then either $i_t \notin Q$ or $i_t \in Q \setminus S_t$. Since S_t contains every $i \in Q$ with $i > j_t$, we can conclude that if $i_t \in Q \setminus S_t$, then $\theta_{i_t} < \theta_{j_t} = p'_t$. Therefore, $\mathbb{I}(\theta_{i_t} \geq p'_t \wedge i_t \in Q \wedge b_t = 0) = 0$, which means that

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t>t_\lambda} p'_t - p_t \right].$$

Let $j'_t = \operatorname{argmin}_{j \in S_t} \operatorname{LB}_{jt}$, which means that $p_t = \min_{i \in S_t} \{\min\{\theta_i, \operatorname{LB}_{it}\}\} = \min\{p'_t, \operatorname{LB}_{j'_t t}\}$. We also know that $p'_t = \theta_{j_t} \leq \theta_{j'_t}$. Therefore,

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t>t_\lambda} \max\{0, \theta_{j'_t} - \operatorname{LB}_{j'_t t}\} \right].$$

Let \mathcal{E}_1 be the event that for all $t > t_\lambda$, $|\Phi_{i,t}| \geq \frac{1}{2}q_{\min}(t-1)$ for all $i \in S_t$. In Lemma B.4, we prove that if $q_{\min} > 2\lambda$, then $\Pr[\mathcal{E}_1^c] \leq \frac{1}{T}$. Also, let \mathcal{H} be the event that for all $t > t_\lambda$ and all $i \in S_t$,

$$|\Phi_{it}| \theta_i \leq \sum_{v \in \Phi_{it}} v + \sqrt{\frac{1}{2} |\Phi_{it}| \ln(dT^2)}.$$

In Lemma B.7, we prove that $\Pr[\mathcal{H}^c] \leq \frac{1}{T}$.

Suppose that $q_{\min} > 2\lambda$. In this case,

$$\begin{aligned} \mathbb{E}[Z_5] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \max\{0, \theta_{j'_t} - \operatorname{LB}_{j'_t t}\} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2 \\ &= \mathbb{E} \left[\sum_{t>t_\lambda} \max \left\{ 0, \theta_{j'_t} - \frac{1}{|\Phi_{j'_t t}|} \sum_{v \in \Phi_{j'_t t}} v + \sqrt{\frac{1}{2|\Phi_{j'_t t}|} \ln \frac{1}{\eta}} \right\} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2. \end{aligned}$$

Under events \mathcal{E}_1 and \mathcal{H} ,

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t>t_\lambda} \sqrt{\frac{\ln(dT^2)}{2|\Phi_{j'_t t}|}} + \sqrt{\frac{1}{2|\Phi_{j'_t t}|} \ln \frac{1}{\eta}} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2$$

and by definition of the event \mathcal{E}_1 ,

$$\begin{aligned} \mathbb{E}[Z_5] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \sqrt{\frac{\ln(dT^2)}{2|\Phi_{j'_t t}|}} + \sqrt{\frac{1}{2|\Phi_{j'_t t}|} \ln \frac{1}{\eta}} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2 \\ &\leq \sum_{t=2}^T \left(\sqrt{\frac{\ln(dT^2)}{q_{\min}(t-1)}} + \sqrt{\frac{1}{q_{\min}(t-1)} \ln \frac{1}{\eta}} \right) + 2 \leq 4\sqrt{\frac{T}{q_{\min}} \ln \frac{dT^2}{\eta}} + 2. \end{aligned}$$

The proof when $q_{\min} < 2\lambda$ follows similarly. □

The final regret bound follows by combining these claims. \square

We conclude this section by providing a proof sketch of one of the lemmas we used in Theorem 4.1. The full proof and the remaining lemmas are in Appendix B. This lemma shows that for all active types $i \in S_t$, $\bar{\mu}_{i,t}$ is indeed a good estimate of $\text{rev}(\theta_i, Q) = \theta_i \Pr_{j \sim \mathcal{P}}[\theta_j \geq \theta_i \text{ and } j \in Q]$. As we described in the proof sketch of Theorem 4.1, this takes some care because we cannot observe at each round t whether or not $i_t \in Q$, provided the buyer did not buy the item.

Lemma 4.7. *For all $t > t_\lambda$, let \mathcal{A}_t be the event $\text{rev}(\theta_i, Q) \in [\check{\mu}_{i,t}, \hat{\mu}_{i,t}]$ for all $i \in S_t$. Then $\Pr[\mathcal{A}_t^c] \leq \frac{1}{T^2}$.*

Proof sketch. Recall that $\hat{\mu}_{i,t} = \bar{\mu}_{i,t} + \rho_t$ and $\check{\mu}_{i,t} = \bar{\mu}_{i,t} - \rho_t$ with

$$\rho_t = \sqrt{\frac{\ln(dT^2)}{2(t - t_\lambda)}}.$$

We also define the related quantities for all $i \in [d]$ and all $Q' \subseteq [d]$:

$$\bar{\gamma}_{i,t}(Q') = \frac{1}{t - t_\lambda} \sum_{s=t_\lambda+1}^t \theta_i \cdot \mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q'),$$

$\hat{\gamma}_{i,t}(Q') = \bar{\gamma}_{i,t}(Q') + \rho_t$, and $\check{\gamma}_{i,t}(Q') = \bar{\gamma}_{i,t}(Q') - \rho_t$. By a Hoeffding bound, for all $Q' \subseteq [d]$ and $i \in [d]$,

$$\Pr[\text{rev}(\theta_i, Q') \notin [\check{\gamma}_{i,t}(Q'), \hat{\gamma}_{i,t}(Q')]] \leq \frac{1}{dT^2}.$$

We claim that for any $i \in S_t$ and any $s > t_\lambda$,

$$\mathbb{I}(b_s = 1 \wedge \theta_{i_s} \geq \theta_i \wedge i_s \in Q) = \mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q), \quad (9)$$

which means that $\bar{\mu}_{i,t} = \bar{\gamma}_{i,t}(Q)$, $\hat{\mu}_{i,t} = \hat{\gamma}_{i,t}(Q)$, and $\check{\mu}_{i,t} = \check{\gamma}_{i,t}(Q)$. To see why, if $b_s = 1$, then clearly Equation (9) holds. Otherwise, suppose $b_s = 0$, in which case $\mathbb{I}(b_s = 1 \wedge \theta_{i_s} \geq \theta_i \wedge i_s \in Q) = 0$. Then $i_s \notin S_s$ because any buyer in S_s will always buy by the definition of the pricing rule. Let $j_s = \min\{j \in S_s\}$. Since S_s contains every element in Q larger than j_s , we know that either:

1. $i_s \notin Q$, in which case $\mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q) = 0$, or
2. $i_s \in Q$ but $i_s \notin S_s$, which means that $\theta_{i_s} < \theta_{j_s}$. Since $i \in S_t$, it must be that $i \in S_s$, so $\theta_{i_s} < \theta_{j_s} \leq \theta_i$. In this case, $\mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q) = 0$ as well.

Therefore, Equation (9) holds.

The fact that $\bar{\mu}_{i,t} = \bar{\gamma}_{i,t}(Q)$, $\hat{\mu}_{i,t} = \hat{\gamma}_{i,t}(Q)$, and $\check{\mu}_{i,t} = \check{\gamma}_{i,t}(Q)$ for all $i \in S_t$ implies that

$$\begin{aligned} \Pr[\mathcal{A}_t^c] &= \Pr(\exists i \in S_t \text{ s.t. } \text{rev}(\theta_i, Q) \notin [\check{\mu}_{i,t}, \hat{\mu}_{i,t}]) \\ &\leq \Pr(\exists i \in [d] \text{ s.t. } \text{rev}(\theta_i, Q) \notin [\check{\gamma}_{i,t}(Q), \hat{\gamma}_{i,t}(Q)]) \\ &\leq \sum_{i=1}^d \Pr(\text{rev}(\theta_i, Q) \notin [\check{\gamma}_{i,t}(Q), \hat{\gamma}_{i,t}(Q)]). \end{aligned} \quad (10)$$

The result now follows from a union bound. \square

5 Regret lower bounds

In this section, we state our regret lower bounds. Recall that $q_{\min} = \min_{i \in [d]} \Pr_{j \sim \mathcal{P}}[j = i]$ denotes the minimum probability of appearance among all types. Let $R_T(A, P)$ denote the regret after T rounds when using an algorithm A on a problem P . Our theorem below presents two lower bounds that correspond to our upper bounds. First, we prove a q_{\min} independent $\tilde{\Omega}(T^{2/3}d^{1/3})$ lower bound on the regret. Next, when q_{\min} is large, we show that $\tilde{\Omega}(\sqrt{T/q_{\min}})$ regret is still unavoidable.

Theorem 5.1. For $T \in \Omega\left(d(\ln(1/\eta))^2(\ln d)^{3/2}\right)$,

$$\inf_A \sup_P R_T(A, P) \geq \frac{1}{4}T^{2/3}(d-1)^{1/3} \left(\ln \frac{1}{\eta}\right)^{1/3} - 2T^{1/2} \in \Omega\left(T^{2/3}d^{1/3} \left(\ln \frac{1}{\eta}\right)^{1/3}\right).$$

Next, suppose $q_{\min} \geq q_T^0 \stackrel{\text{def}}{=} T^{-1/3}(d-1)^{-2/3}(\ln(1/\eta))^{1/3}$. Then for $T \in \Omega\left(d(\ln(1/\eta))^2(\ln d)^{3/2}\right)$,

$$\inf_A \sup_{P; q_{\min} \geq q_T^0} R_T(A, P) \geq \frac{1}{4}\sqrt{\frac{T}{q_{\min}} \ln \frac{1}{\eta}} - 2T^{1/2} \in \Omega\left(\sqrt{\frac{T}{q_{\min}} \ln \frac{1}{\eta}}\right).$$

Comparing this with Corollary 4.2, we see that our algorithm is minimax optimal, up to constants and polylog terms. This is the case even when q_{\min} is larger than $\tilde{\Omega}(T^{-1/3}d^{-2/3})$ where \sqrt{T} rates are possible. As we mentioned at the end of Section 1.1, our proof reveals interesting properties about the structure of an optimal policy; we discuss these in detail at the end of this section.

Proof of Theorem 5.1. Unlike typical proofs of lower bounds in stochastic bandit settings, which usually rely on hypothesis testing arguments, our result stems from the buyers' uncertainty about their values. To demonstrate this, we will construct a representative problem instance and show that any algorithm will do poorly on this instance.

Construction. For all types $j \in [d]$, we set the *ex-post* value distribution to be $\mathcal{D}_j = \text{Unif}(1 - 2/\sqrt{T}, 1)$. Hence, for all $j \in [d]$, $\theta_j = 1 - 1/\sqrt{T}$. Next, we define the type distribution \mathcal{P} as shown below. Here $q < 1/d$ is a parameter we will specify later in the proof.

$$\forall j \in \{1, \dots, d-1\}, q_j = \Pr_{i \sim \mathcal{P}}[i = j] = q, \quad q_d = \Pr_{i \sim \mathcal{P}}[i = d] = 1 - q(d-1). \quad (11)$$

We will use the following threshold functions for each buyer of each type. Recall that $\Phi_{i,t-1} = \{v; (i, v) \in \sigma_{t-1}\}$ denotes the reviews in σ_{t-1} left by customers of type i .

$$\tau_t(\sigma_{t-1}, i) = \max \left\{ \frac{1}{|\Phi_{i,t-1}|} \sum_{v \in \Phi_{i,t-1}} v - \sqrt{\frac{1}{2|\Phi_{i,t-1}|} \ln \frac{1}{\eta}}, 0 \right\}.$$

Note that $\tau_t(\sigma_{t-1}, i_t)$ is larger than LB_t as defined in Definition 2.1 and satisfies the η -pessimistic agents' assumption. We will also assume that the seller knows the type distribution \mathcal{P} ; this additional information can only help the seller. Despite this, we show that if buyers choose conservative threshold functions, $T^{2/3}$ regret is unavoidable.

The optimal price for the above construction is $p^* = \theta_1 = \dots = \theta_d = 1 - 1/\sqrt{T}$. The seller could simply set this price if all buyers knew their *ex-ante* values. However, when buyers learn their values from past observations, the confidence of their estimates shrinks only with the number of observations of their type. In particular, if q is very small, then a seller might find it beneficial to ignore customers of the first $d - 1$ types and set the highest possible price that can still attract customers of type d . On the other hand, if q is large, the higher price may not warrant the revenue foregone by ignoring the first $d - 1$ types. By carefully choosing q , we can balance these trade-offs to obtain the tightest lower bound. We have set the value of all types to be equal in this construction to simplify some of our calculations, but it is not hard to see how this phenomenon affects pricing decisions for the seller.

Set up and notation. For brevity, we use the following notation for the sample mean of observations, the number of observations, and the threshold function for type i on round t .

$$\begin{aligned} \hat{\mathbb{E}}_{i,t-1} &\stackrel{\text{def}}{=} \frac{1}{|\Phi_{i,t-1}|} \sum_{v \in \Phi_{i,t-1}} v, & \tilde{N}_{i,t-1} &\stackrel{\text{def}}{=} |\Phi_{i,t-1}|, \\ \tau_{i,t} &\stackrel{\text{def}}{=} \tau_t(\sigma_{t-1}, i) = \max \left\{ \hat{\mathbb{E}}_{i,t-1} - \sqrt{\frac{1}{2\tilde{N}_{i,t-1}} \ln \left(\frac{1}{\eta} \right)}, 0 \right\}. \end{aligned} \quad (12)$$

Next, let τ'_t denote the maximum of the threshold functions of the first $d - 1$ types on round t and i'_t denote the corresponding maximizer.

$$i'_t = \operatorname{argmax}_{j \in \{1, \dots, d-1\}} \tau_{j,t}, \quad \tau'_t = \tau_{i'_t,t}. \quad (13)$$

Recall that on each round t , a seller's policy chooses a price p_t based on all past information σ_{t-1} and possibly some source of external randomness. We next define $W_{1,t}, W_{2,t}, W_{3,t}$ below based on how p_t compares to the threshold functions:

$$W_{1,t} = \mathbb{I}(p_t \leq \tau'_t), \quad W_{2,t} = \mathbb{I}(\tau'_t < p_t \leq \tau_{d,t}), \quad W_{3,t} = \mathbb{I}(\forall j \in [d], \tau_{j,t} < p_t). \quad (14)$$

Here, $W_{1,t}$ is 1 when the price p_t is smaller than the thresholds for any of the first $d - 1$ types, $W_{2,t}$ is 1 when p_t is larger than the thresholds for all $d - 1$ types but smaller than the threshold $\tau_{d,t}$ for type d (note that $W_{2,t}$ can be 1 only when $\tau'_t < \tau_{d,t}$), and $W_{3,t}$ is 1 when p_t is larger than all thresholds. It is easy to verify that exactly one of $W_{1,t}, W_{2,t}, W_{3,t}$ is 1 on any given round.

Lower bounding the instantaneous regret. We can decompose the expected revenue $\operatorname{rev}_t = p_t b_t$ on round t , conditioned on the price p_t and history σ_{t-1} as follows.

$$\begin{aligned} \mathbb{E}[\operatorname{rev}_t | \sigma_{t-1}, p_t] &= \mathbb{E}[p_t b_t | \sigma_{t-1}, p_t] \\ &= \sum_{i \in [d]} p_t \mathbb{E}[b_t | \sigma_{t-1}, p_t, i_t = i] \Pr[i_t = i | \sigma_{t-1}, p_t] \\ &= \sum_{i \in [d]} p_t \mathbb{I}(p_t \leq \tau_{i,t}) q_i \\ &= \sum_{i \in \{(d-1)\}} p_t \mathbb{I}(p_t \leq \tau_{i,t}) q + p_t \mathbb{I}(p_t \leq \tau_{d,t}) q_d. \end{aligned} \quad (15)$$

In the third step we have used the fact that the probability of appearance of a type does not depend on the history or the price chosen, hence $\Pr[i_t = i | \sigma_{t-1}, p_t] = \Pr_{j \sim \mathcal{P}}[j = i] = q_i$. Second, we note that for a customer of type i , they will purchase if and only if the price is smaller than their threshold; therefore $\mathbb{E}[b_t | \sigma_{t-1}, p_t, i_t = i] = \mathbb{1}(p_t \leq \tau_{i,t})$. The following lemma upper bounds $\mathbb{E}[\text{rev}_t | p_t]$ in terms of the $W_{i,t}$ terms defined in (14).

Lemma 5.2. $\mathbb{E}[\text{rev}_t | \sigma_{t-1}, p_t] \leq W_{1,t} \tau'_t + W_{2,t} q_d \tau_{d,t}$.

Proof of Lemma 5.2. We will consider four exhaustive cases for p_t and analyze the right-hand side of the inequality in the claim as a function of $W_{1,t}$ and $W_{2,t}$, which we denote as $\text{RHS}(W_{1,t}, W_{2,t})$.

1. $p_t \leq \min\{\tau'_t, \tau_{d,t}\}$: Here, $W_{1,t} = 1$ and $W_{2,t} = 0$. Using (15), we obtain $\mathbb{E}[\text{rev}_t | p_t] \leq (d-1)p_t q + p_t q_d = p_t \leq \tau'_t = \text{RHS}(1, 0)$.
2. $\tau'_t < p_t \leq \tau_{d,t}$: Here, $W_{1,t} = 0$ and $W_{2,t} = 1$. Using (15), $\mathbb{E}[\text{rev}_t | p_t] = p_t q_d \leq \tau_{d,t} q_d = \text{RHS}(0, 1)$.
3. $\tau_{d,t} < p_t \leq \tau'_t$: Here, $W_{1,t} = 1$ and $W_{2,t} = 0$. Using (15), we obtain

$$\mathbb{E}[\text{rev}_t | p_t] \leq p_t (d-1)q < p_t \leq \tau'_t = \text{RHS}(1, 0).$$

Some of terms in the first summation in (15) may be 0, but we can bound it by $p_t (d-1)q$ regardless.

4. $p_t > \max\{\tau'_t, \tau_{d,t}\}$: Here, $W_{1,t} = W_{2,t} = 0$. Using (15), $\mathbb{E}[\text{rev}_t | p_t] = 0 = \text{RHS}(0, 0)$.

□

Equipped with this lemma, we can now lower bound the instantaneous regret on round t conditioned on the price p_t and history σ_{t-1} , which we denote as $\mathbb{E}[r_t | \sigma_{t-1}, p_t]$:

$$\begin{aligned} \mathbb{E}[r_t | \sigma_{t-1}, p_t] &= \mathbb{E}[p^* - \text{rev}_t | \sigma_{t-1}, p_t] \\ &\geq W_{1,t} \cdot (p^* - \tau'_t) + W_{2,t} \cdot (p^* (d-1)q + q_d (p^* - \tau_{d,t})) + W_{3,t} \cdot p^*. \end{aligned} \quad (16)$$

Recall that i'_t is the index such that $\tau'_t = \tau_{i'_t,t}$ as defined in (13) and $\tilde{N}_{i'_t,t-1}$ is the number of observations of type i'_t in σ_{t-1} as defined in (12). We can further lower bound Equation (16) by using the fact that

$$p^* - \tau'_t = p^* - \tau_{i'_t,t} = p^* - \max \left\{ \hat{\mathbb{E}}_{i'_t,t-1} - \sqrt{\frac{1}{2\tilde{N}_{i'_t,t-1}} \ln \left(\frac{1}{\eta} \right)}, 0 \right\}. \quad (17)$$

Since the support of each *ex-post* value distribution is bounded within an $\pm 1/\sqrt{T}$ interval of p^* , Equation (17) implies that

$$\begin{aligned} p^* - \tau'_t &\geq p^* - \max \left\{ p^* + \frac{1}{\sqrt{T}} - \sqrt{\frac{1}{2\tilde{N}_{i'_t,t-1}} \ln \left(\frac{1}{\eta} \right)}, 0 \right\} \\ &= \min \left\{ \sqrt{\frac{1}{2\tilde{N}_{i'_t,t-1}} \ln \left(\frac{1}{\eta} \right)} - \frac{1}{\sqrt{T}}, p^* \right\} \\ &= \min \left\{ \sqrt{\frac{1}{2\tilde{N}_{i'_t,t-1}} \ln \left(\frac{1}{\eta} \right)}, 1 \right\} - \frac{1}{\sqrt{T}}. \end{aligned} \quad (18)$$

The same argument guarantees that $p^* - \tau_{d,t} \geq -\frac{1}{\sqrt{T}}$. Combining this inequality with Equations (16) and (17), and recalling that $W_{1,t} + W_{2,t} + W_{3,t} = 1$, we have that

$$\mathbb{E}[r_t | \sigma_{t-1}, p_t] \geq W_{1,t} \min \left\{ \sqrt{\frac{1}{2\tilde{N}_{i'_t, t-1}}} \ln \left(\frac{1}{\eta} \right), 1 \right\} + W_{2,t} p^* (d-1)q + W_{3,t} p^* - \frac{1}{\sqrt{T}}. \quad (19)$$

Upper bounding $\tilde{N}_{i'_t, t-1}$. To convert the above instantaneous bound to a lower bound on the cumulative regret, we will need to control $\tilde{N}_{i'_t, t-1}$ which counts the number of reviews in σ_{t-1} by customers of type i'_t . Observing that $i'_t \in [(d-1)]$ which means that the appearance probability of i'_t is q , we define the following event \mathcal{E} below. Lemma 5.3 upper bounds the probability of this event.

$$\mathcal{E} = \left\{ \forall j \in [(d-1)], \forall t \leq T, \tilde{N}_{j, t-1} \leq 2q(T-1) \right\}. \quad (20)$$

Lemma 5.3. *Let $T \geq \frac{3}{q} \ln(2d) + 1$. Then, $\Pr[\mathcal{E}] \geq 1/2$.*

Proof of Lemma 5.3. Note that

$$\tilde{N}_{i, t-1} = \sum_{s=1}^{t-1} \mathbb{1}(b_s = 1, i_s = i)$$

counts the number of times a customer of type i made a purchase. Let

$$N_{i, t-1} = \sum_{s=1}^{t-1} \mathbb{1}(i_s = i)$$

be the number of times a customer of type i arrived. Since $N_{i, T-1} \geq \tilde{N}_{i, T-1}$, the Chernoff bound implies that

$$\Pr[\exists t \leq T \text{ such that } \tilde{N}_{i, t-1} > 2q(T-1)] \leq \Pr[N_{i, T-1} > 2q(T-1)] \leq \exp\left(-\frac{q(T-1)}{3}\right) \leq \frac{1}{2d}.$$

The last step uses the condition on T . The claim follows via a union bound over $j \in [d-1]$. \square

Lower bound on cumulative regret. We are now ready to lower bound regret. By Equation (19),

$$\mathbb{E}[R_T] \geq \mathbb{E} \left[\sum_{t=1}^T \left(\underbrace{W_{1,t} \min \left\{ \sqrt{\frac{1}{2\tilde{N}_{i'_t, t-1}}} \ln \left(\frac{1}{\eta} \right), 1 \right\}}_{\bar{r}_t} + W_{2,t} p^* (d-1)q + W_{3,t} p^* - \frac{1}{\sqrt{T}} \right) \right].$$

Conditioning on the event \mathcal{E} ,

$$\mathbb{E}[R_T] \geq -\sqrt{T} + \mathbb{E} \left[\sum_{t=1}^T \bar{r}_t \mid \mathcal{E} \right] \Pr[\mathcal{E}] + \mathbb{E} \left[\sum_{t=1}^T \bar{r}_t \mid \mathcal{E}^c \right] \Pr[\mathcal{E}^c]$$

and by Lemma 5.3,

$$\begin{aligned}\mathbb{E}[R_T] &\geq -\sqrt{T} + \frac{1}{2} \mathbb{E} \left[\sum_{t=1}^T \bar{r}_t \mid \mathcal{E} \right] - T \cdot \frac{1}{\sqrt{T}} \\ &\geq -2\sqrt{T} + \frac{1}{2} \mathbb{E} \left[\sum_{t=1}^T W_{1,t} \min \left\{ \sqrt{\frac{1}{4qT} \ln \left(\frac{1}{\eta} \right)}, 1 \right\} + W_{2,t} p^*(d-1)q + W_{3,t} p^* \mid \mathcal{E} \right].\end{aligned}$$

For $T > \frac{1}{4q} \ln \left(\frac{1}{\eta} \right)$,

$$\mathbb{E}[R_T] \geq -2\sqrt{T} + \frac{1}{2} \mathbb{E} \left[\sum_{t=1}^T W_{1,t} \sqrt{\frac{1}{4qT} \ln \left(\frac{1}{\eta} \right)} + W_{2,t} p^*(d-1)q + W_{3,t} p^* \mid \mathcal{E} \right].$$

We will use the notation $M_{i,t} = \sum_{s=1}^t W_{i,s}$ for $i \in \{1, 2, 3\}$ which counts the number of times each $W_{i,s}$ was 1 in the first t rounds. Note that $M_{1,t} + M_{2,t} + M_{3,t} = t$ since exactly one of $W_{1,s}, W_{2,s}, W_{3,s}$ is 1 on any round s . With this notation, we have that

$$\mathbb{E}[R_T] \geq -2\sqrt{T} + \frac{1}{2} \mathbb{E} \left[M_{1,T} \sqrt{\frac{1}{4qT} \ln \left(\frac{1}{\eta} \right)} + M_{2,T} p^*(d-1)q + M_{3,T} p^* \mid \mathcal{E} \right]. \quad (21)$$

We note that $M_{1,T}, M_{2,T}, M_{3,T}$ are random quantities that depend on the execution of the algorithm. However, we can use the fact that they are non-negative and that $M_{1,T} + M_{2,T} + M_{3,T} = T$ to obtain a lower bound as follows.

$$\mathbb{E}[R_T] \geq -2\sqrt{T} + \frac{1}{2} \inf_{\substack{x_1, x_2, x_3 > 0 \\ x_1 + x_2 + x_3 = T}} \left(x_1 \sqrt{\frac{1}{4qT} \ln \left(\frac{1}{\eta} \right)} + x_2 p^*(d-1)q + x_3 p^* \right).$$

As $(d-1)q \leq 1$, for any choice (x'_1, x'_2, x'_3) for (x_1, x_2, x_3) such that $x'_3 > 0$, we can obtain a lower value for the term in parentheses via $(x'_1, x'_2 + x'_3, 0)$. Therefore, the above expression simplifies to:

$$\mathbb{E}[R_T] \geq -2\sqrt{T} + \frac{1}{2} \inf_{0 \leq x \leq T} \left(x \sqrt{\frac{1}{4qT} \ln \left(\frac{1}{\eta} \right)} + (T-x) p^*(d-1)q \right). \quad (22)$$

Finally, we are taking the infimum of a linear function in the bounded interval $[0, T]$, so the infimum lies at one of the end points $x = 0$ or $x = T$. Therefore,

$$\begin{aligned}\mathbb{E}[R_T] &\geq -2\sqrt{T} + \frac{1}{2} \min \left\{ \sqrt{\frac{T}{4q} \ln \left(\frac{1}{\eta} \right)}, T p^*(d-1)q \right\} \\ &\geq -2\sqrt{T} + \frac{1}{2} \min \left\{ \sqrt{\frac{T}{4q} \ln \left(\frac{1}{\eta} \right)}, T p^*(d-1)q \right\}.\end{aligned} \quad (23)$$

Putting it all together. To complete the proof, first note that for all $T \geq 4$, $p^* \geq 1/2$; hence, the second term inside the min can be upper bounded by $\frac{1}{2} T (d-1)q$. To obtain a q_{\min} independent bound, we set $q = T^{-1/3} (d-1)^{-2/3} (\ln(1/\eta))^{1/3}$ to obtain the first result of the theorem.

Next, since $q_{\min} = q$ for this problem, we have that when

$$q_{\min} > q_T^0 = T^{-1/3}(d-1)^{-2/3}(\ln(1/\eta))^{1/3},$$

the minimum is the first of the two terms in (23). This leads to our second lower bound. \square

Our construction uses p^* close to 1 to simplify some of the calculations in the analysis, but a similar analysis is possible for any p^* bounded away from 0. Second, while our construction sets the ex-ante value θ_j to be the same for all types, a similar result can be shown in cases where a low probability type has ex-ante value similar to or larger than the ex-ante value of high probability types. Third, recall that we have assumed in this proof that the seller knows the type distribution \mathcal{P} . If it is unknown, as was shown in our upper-bound analysis, the seller only really needs to estimate the low probability types and the expected revenue when targeting the remaining types, both of which can be done at rates $T^{1/3}$ and $T^{1/2}$ respectively without having to learn \mathcal{P} entirely. The $T^{2/3}$ bottleneck arises as the seller needs to wait for the buyers' estimates of their values become accurate.

We also make the following observation via Equations (21)–(23). Intuitively, $M_{1,T}$ in (21) denotes the number of times the seller's policy targeted the low probability types, $M_{2,T}$ denotes the number of times it targeted the high probability type while ignoring the low probability types, and $M_{3,T}$ is the number of times it targeted none of the types. Equation (22) states that any reasonable policy will never ignore all customer types, choosing $M_{3,T} = 0$. On the other hand, the fact that the infimum in (23) lies in one of two extremes $(M_{1,T}, M_{2,T}) \in \{(0, T), (T, 0)\}$ indicates that any reasonable policy cannot do significantly better than a policy which chooses ahead of time to target all customer types or only focus on the high probability types. Intuitively, this means that the seller's policy can decide ahead of time which customers it wants to ignore due to a low probability of appearance. In other words, it does not significantly help to change which types you target on different rounds based on their appearance probability. Interestingly, this is precisely the behavior of our algorithm as well; it uses a small initial phase of at most $T^{1/3}$ rounds to identify and eliminate low probability types. From thereon, it only targets the remaining high probability types.

6 Conclusion

We proposed no-regret online pricing strategies when both sides of the market learn from reviews. Our algorithm strategically sets lower prices during its early phase to boost sales from customers with rare types and high values. Reviews from the early phase benefit future buyers in the long run. Our algorithm carefully trades off the revenue loss due to discounts from the initial phase and future gains. Our lower bound demonstrates that our algorithm is optimal up to lower order and constant terms. To the best of our knowledge, this is the first result on online pricing when both the seller learns to price and buyers with different types learn from reviews.

Future directions. Many questions remain open for future research. We assumed that purchases always come with a noisy review. An interesting direction would be providing pricing strategies when the reviews are left with varying probabilities, which mimics real-world buyer behaviors.

We studied myopic buyers who make their purchase decisions based on estimates of their *ex-ante* values from historical reviews, regardless of the seller's policy. What if the buyers appear over several rounds and may behave strategically to purchase at lower future prices?

We take a frequentist perspective on this problem. It is also possible to take a Bayesian view of this problem and impose a prior on the *ex-ante* value so that the buyer starts with some prior information. We expect adapting our main proof intuitions to that setting is possible. The main differences would be: (i) we would use Bayesian credible intervals instead of frequentist confidence intervals for the η -pessimism definition, (ii) we would control the Bayes' risk when estimating the *ex-ante* values instead of frequentist concentration arguments, and (iii) our final regret could have a nuanced dependence on this prior which may offer tighter bounds.

Another direction would be to explore the case where the seller does not know the buyers' *ex-ante* values. The key challenge would be related to the regret benchmark: we compete with the optimal price if the buyers knew their own *ex-ante* values and bought whenever their *ex-ante* value was above the price (thus, the buyers are not learning). To compete with this benchmark, we require unbiased estimates of the revenue of different prices if the buyers bought when their *ex-ante* value was above the price. Computing these unbiased estimates is challenging: if a buyer does not buy on a given round, the algorithm does not learn their type, so it cannot tell whether the buyer has a low *ex-ante* value or he has a high value but a low confidence bound. If the seller knows the buyers' *ex-ante* values, we can circumvent this subtle challenge, as we explain in the proof sketch of Theorem 4.1. However, this is not possible if the *ex-ante* values are unknown.

References

- Daron Acemoglu, Ali Makhdoumi, Azarakhsh Malekian, and Asuman Ozdaglar. Learning from reviews: The selection effect and the speed of learning. *Econometrica*, 90(6):2857–2899, 2022.
- Itai Ashlagi, Constantinos Daskalakis, and Nima Haghpanah. Sequential mechanisms with ex-post participation guarantees. In *ACM Conference on Economics and Computation (EC)*, 2016.
- Omar Besbes and Marco Scarsini. On information distortions in online ratings. *Operations Research*, 66(3):597–610, 2018.
- Subir Bose, Gerhard Orosel, Marco Ottaviani, and Lise Vesterlund. Dynamic monopoly pricing and herding. *The RAND Journal of Economics*, 37(4):910–928, 2006.
- Etienne Boursier, Vianney Perchet, and Marco Scarsini. Social learning in non-stationary environments. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 128–129, 2022.
- Mark Braverman, Jieming Mao, Jon Schneider, and Matt Weinberg. Selling to a no-regret buyer. In *ACM Conference on Economics and Computation (EC)*, 2018.
- Christophe Chamley. *Rational herds: Economic models of social learning*. Cambridge University Press, 2004.
- Shuchi Chawla, Nikhil R Devanur, Anna R Karlin, and Balasubramanian Sivan. Simple pricing schemes for consumers with evolving values. *Games and Economic Behavior*, 134:344–360, 2022.
- Davide Crippa, Bar Ifrach, Costis Maglaras, and Marco Scarsini. Monopoly pricing in the presence of social learning. *Management Science*, 63(11):3586–3608, 2017.

- Yuan Deng, Jon Schneider, and Balasubramanian Sivan. Prior-free dynamic auctions with low regret buyers. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- Nikhil R Devanur, Yuval Peres, and Balasubramanian Sivan. Perfect Bayesian equilibria in repeated sales. In *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2014.
- Zhe Feng, Chara Podimata, and Vasilis Syrgkanis. Learning to bid without knowing your value. In *ACM Conference on Economics and Computation (EC)*, 2018.
- Aurélien Garivier, Tor Lattimore, and Emilie Kaufmann. On explore-then-commit strategies. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2016.
- Saram Han and Chris K Anderson. Customer motivation and response bias in online reviews. *Cornell Hospitality Quarterly*, 61(2):142–153, 2020.
- Bar Ifrach, Costis Maglaras, Marco Scarsini, and Anna Zseleva. Bayesian social learning from consumer reviews. *Operations Research*, 67(5):1209–1221, 2019.
- Ali Kakhbod, Giacomo Lanzani, and Hao Xing. Heterogeneous Learning in Product Markets. *Available at SSRN 3961223*, 2021.
- Kirthevasan Kandasamy, Joseph E Gonzalez, Michael I Jordan, and Ion Stoica. VCG mechanism design with unknown agent values under stochastic bandit feedback. *Journal of Machine Learning Research*, 24(53):1–45, 2023.
- Robert Kleinberg and Tom Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *Symposium on Foundations of Computer Science (FOCS)*, 2003.
- Denis Nekipelov, Vasilis Syrgkanis, and Eva Tardos. Econometrics for learning agents. In *ACM Conference on Economics and Computation (EC)*, 2015.
- Christos Papadimitriou, George Pierrakos, Alexandros Psomas, and Aviad Rubinfeld. On the complexity of dynamic mechanism design. *Games and Economic Behavior*, 2022.
- Jonathan Weed, Vianney Perchet, and Philippe Rigollet. Online learning in repeated auctions. In *Conference on Learning Theory (COLT)*, 2016.
- Haoyu Zhao and Wei Chen. Stochastic one-sided full-information bandit. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, 2020.

A Additional details about η -pessimistic agents

Intuitively, in Definition 2.1, LB_t serves as a lower confidence bound on the buyer’s value who arrives at round t . The buyers can be confident that, *regardless* of the policy used by the seller, with probability $1 - \eta$, for all rounds $t \in [T]$, $\theta_{i_t} \geq \text{LB}_t$. We show this formally below.

Lemma A.1. *Denote the type of the buyer who arrives at round t as i_t . On all rounds t , with probability at least $1 - \eta$, $\text{LB}_t \leq \theta_{i_t}$.*

Proof. Let us consider a sequence of T rewards $\{\tilde{v}_{i1}, \dots, \tilde{v}_{iT}\}$ for each buyer type $i \in [d]$ generated beforehand, where each reward is a random reward sample drawn from \mathcal{D}_i . Each time a buyer with type i arrives and makes a purchase, it obtains an ex-post value from the reward sequence $\{\tilde{v}_{i1}, \dots, \tilde{v}_{iT}\}$ in order. For example, if the type of the buyer who arrives on round t is i_t , then if that buyer makes a purchase, their ex-post value will be $\tilde{v}_{i_t, |\Phi_t|+1}$.

First, we will show that $\Pr(\text{LB}_t > \theta_j \mid i_t = j) \leq \eta$ for any $j \in [d]$. At any round t , notice that if $|\Phi_t| = 0$, then $\text{LB}_t = 0$, the conclusion trivially holds since $\theta_j > 0$ for all $j \in [d]$. When $|\Phi_t| > 0$:

$$\begin{aligned}
& \Pr\left(\text{LB}_t > \theta_j \mid i_t = j\right) \\
&= \Pr\left(\max\left\{0, \frac{1}{|\Phi_t|} \sum_{v \in \Phi_t} v - \sqrt{\frac{1}{2|\Phi_t|} \ln \frac{t}{\eta}}\right\} > \theta_j \mid i_t = j\right) \\
&= \Pr\left(\frac{1}{|\Phi_t|} \sum_{v \in \Phi_t} v - \sqrt{\frac{1}{2|\Phi_t|} \ln \frac{t}{\eta}} > \theta_j \mid i_t = j\right) \\
&= \Pr\left(\frac{1}{|\Phi_t|} \sum_{s=1}^{|\Phi_t|} \tilde{v}_{js} - \sqrt{\frac{1}{2|\Phi_t|} \ln \frac{t}{\eta}} > \theta_j \mid i_t = j\right) \\
&\leq \Pr\left(\exists \ell \in [t-1], \text{s.t. } \frac{1}{\ell} \sum_{s=1}^{\ell} \tilde{v}_{js} - \sqrt{\frac{1}{2\ell} \ln \frac{t}{\eta}} > \theta_j \mid i_t = j\right) \\
&\leq \sum_{\ell=1}^{t-1} \Pr\left(\frac{1}{\ell} \sum_{s=1}^{\ell} \tilde{v}_{js} - \sqrt{\frac{1}{2\ell} \ln \frac{t}{\eta}} > \theta_j \mid i_t = j\right).
\end{aligned}$$

Here, the second step uses the fact that $\theta_j \geq 0$. In the fifth step, we have used the fact that $|\Phi_t|$ is a random quantity, which depends on the specific algorithm, but with support $[(t-1)]$. The last step follows from a union bound over $(t-1)$ rounds.

Note that for any fixed $j \in [d]$, the event $\frac{1}{\ell} \sum_{s=1}^{\ell} \tilde{v}_{js} - \sqrt{\frac{1}{2\ell} \ln \frac{t}{\eta}} > \theta_j$ is independent of the value of i_t . Therefore, by Hoeffding inequality, for any $\ell \in [t-1]$ and $j \in [d]$, we have that

$$\Pr\left(\frac{1}{\ell} \sum_{s=1}^{\ell} \tilde{v}_{js} - \sqrt{\frac{1}{2\ell} \ln \frac{t}{\eta}} > \theta_j \mid i_t = j\right) = \Pr\left(\frac{1}{\ell} \sum_{s=1}^{\ell} \tilde{v}_{js} - \sqrt{\frac{1}{2\ell} \ln \frac{t}{\eta}} > \theta_j\right) \leq \frac{\eta}{t}.$$

Putting this together we have:

$$\Pr(\text{LB}_t > \theta_j \mid i_t = j) \leq (t-1) \frac{\eta}{t} \leq \eta.$$

Lastly, by the law of total probability,

$$\Pr(\text{LB}_t > \theta_{i_t}) = \sum_{j \in [d]} \Pr(\text{LB}_t > \theta_j \mid i_t = j) \cdot \Pr(i_t = j) \leq \sum_{j \in [d]} \eta \cdot \Pr(i_t = j) \leq \eta,$$

which completes the proof. \square

B Additional proofs about regret upper bound (Section 4.1)

Claim 4.3. *If $q_{\min} \leq 2\lambda$ then $\mathbb{E}[Z_2] \leq Td\lambda + 1$ and if $q_{\min} > 2\lambda$, then $\mathbb{E}[Z_2] \leq 1$.*

Proof. First, we bound Z_2 as follows:

$$Z_2 = \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \notin Q) \leq \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(i_t \notin Q) \leq \sum_{t>t_\lambda} \mathbb{I}(i_t \notin Q).$$

Recall from Algorithm 1 that \bar{q}_i is the fraction of times that type i appears in phase 1 and let \mathcal{G} be the event that for all $i \in [d]$ such that $q_i \geq \lambda$, we have that $\bar{q}_i \geq \frac{3\lambda}{4}$, which means that $i \in Q$. In other words, when \mathcal{G} happens, $[d] \setminus Q \subseteq \{q_i : q_i < \lambda\}$. In Lemma B.1, we prove that $\Pr[\mathcal{G}^c] \leq \frac{1}{T}$, so $\mathbb{E}[Z_2] \leq \mathbb{E}[Z_2 \mid \mathcal{G}] + T \Pr[\mathcal{G}^c] \leq \mathbb{E}[Z_2 \mid \mathcal{G}] + 1$.

Next, since Q is a random variable, we condition on it as well:

$$\mathbb{E}[Z_2 \mid \mathcal{G}] = \sum_{Q' \subseteq [d]} \mathbb{E}[Z_2 \mid Q = Q', \mathcal{G}] \Pr[Q = Q' \mid \mathcal{G}].$$

If $[d] \setminus Q' \not\subseteq \{q_i : q_i < \lambda\}$, then $\Pr[Q = Q' \mid \mathcal{G}] = 0$. Moreover, for any Q' such that $[d] \setminus Q' \subseteq \{q_i : q_i < \lambda\}$,

$$\begin{aligned} \mathbb{E}[Z_2 \mid Q = Q', \mathcal{G}] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \mathbb{I}(i_t \notin Q') \mid Q = Q', \mathcal{G} \right] \\ &= \sum_{t>t_\lambda} \Pr[i_t \notin Q' \mid Q = Q', \mathcal{G}] \\ &= \sum_{t>t_\lambda} \sum_{i \notin Q'} \Pr[i_t = i \mid Q = Q', \mathcal{G}]. \end{aligned}$$

The event $(Q = Q' \wedge \mathcal{G})$ depends only on the first t_λ timesteps, so it is independent of the event that $i_t = i$ for $t > t_\lambda$. Therefore,

$$\mathbb{E}[Z_2 \mid Q = Q', \mathcal{G}] = \sum_{t>t_\lambda} \sum_{i \notin Q'} \Pr[i_t = i].$$

If $q_{\min} > 2\lambda$, then $\{q_i : q_i < \lambda\} = \emptyset$, so the only Q' such that $[d] \setminus Q' \subseteq \{q_i : q_i < \lambda\}$ is $Q' = [d]$. In this case,

$$\sum_{t>t_\lambda} \sum_{i \notin Q'} \Pr[i_t = i] = 0,$$

so $\mathbb{E}[Z_2 \mid \mathcal{G}] = 0$ and finally, $\mathbb{E}[Z_2] \leq 1$.

Otherwise, $q_{\min} \leq 2\lambda$, so

$$\sum_{t>t_\lambda} \sum_{i \notin Q'} \Pr[i_t = i] \leq \sum_{t>t_\lambda} \sum_{i \notin Q'} \lambda \leq Td\lambda,$$

so

$$\mathbb{E}[Z_2 \mid \mathcal{G}] \leq Td\lambda \sum_{Q' \subseteq [d]} \Pr[Q = Q' \mid \mathcal{G}] \leq Td\lambda$$

and finally, $\mathbb{E}[Z_2] \leq Td\lambda + 1$. □

Claim 4.4. $\mathbb{E}[Z_3] \leq 0$.

Proof. In this proof we bound

$$Z_3 = \sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q). \quad (24)$$

We begin by conditioning the first term of Equation (24) on Q since it is a random variable:

$$\begin{aligned} & \mathbb{E} \left[\sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) \right] \\ &= \sum_{Q' \subseteq [d]} p^*([d]) \mathbb{E} \left[\sum_{t>t_\lambda} \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q') \mid Q = Q' \right] \Pr[Q = Q'] \\ &= \sum_{Q' \subseteq [d]} \sum_{t>t_\lambda} p^*([d]) \Pr[\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q' \mid Q = Q'] \Pr[Q = Q']. \end{aligned} \quad (25)$$

The event that $Q = Q'$ only depends on the first t_λ timesteps, so it is independent of the event $(\theta_{i_t} \geq p^*([d]) \wedge i_t \in Q')$ for $t > t_\lambda$. Therefore, for $t > t_\lambda$,

$$\begin{aligned} p^*([d]) \Pr[\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q' \mid Q = Q'] &= p^*([d]) \Pr[\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q'] \\ &\leq \max_{p \in [0,1]} p \Pr[\theta_{i_t} \geq p \text{ and } i_t \in Q'] \\ &= p^*(Q') \Pr[\theta_{i_t} \geq p^*(Q') \text{ and } i_t \in Q']. \end{aligned}$$

Combining this fact with Equation (25), we have that

$$\begin{aligned} & \mathbb{E} \left[\sum_{t>t_\lambda} p^*([d]) \mathbb{I}(\theta_{i_t} \geq p^*([d]) \text{ and } i_t \in Q) \right] \\ &\leq \sum_{Q' \subseteq [d]} (T - t_\lambda) p^*(Q') \Pr_{i \sim \mathcal{P}}[\theta_i \geq p^*(Q') \text{ and } i \in Q'] \Pr[Q = Q']. \end{aligned} \quad (26)$$

Next, for the second term of Equation (24),

$$\begin{aligned} & \mathbb{E} \left[\sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q) \right] \\ &= \sum_{Q' \subseteq [d]} \mathbb{E} \left[\sum_{t>t_\lambda} p^*(Q') \mathbb{I}(\theta_{i_t} \geq p^*(Q') \text{ and } i_t \in Q') \mid Q = Q' \right] \Pr[Q = Q']. \end{aligned}$$

As before, the event that $Q = Q'$ only depends on the first t_λ timesteps, so it is independent of the event $(\theta_{i_t} \geq p^*(Q') \wedge i_t \in Q')$ for $t > t_\lambda$. Therefore,

$$\begin{aligned} & \mathbb{E} \left[\sum_{t>t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q) \right] \\ &= \sum_{Q' \subseteq [d]} (T - t_\lambda) p^*(Q') \Pr_{i \sim \mathcal{P}}[\theta_i \geq p^*(Q') \text{ and } i \in Q'] \Pr[Q = Q']. \end{aligned}$$

Combined with Equation (26), we have that $\mathbb{E}[Z_3] \leq 0$. □

Claim 4.5. $\mathbb{E}[Z_4] \leq 5 + 4\sqrt{2T \ln(dT^2)}$.

Proof. In this claim, we bound

$$Z_4 = \sum_{t > t_\lambda} p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q) - \sum_{t > t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q), \quad (27)$$

where $p'_t = \min_{i \in S_t} \theta_i$. Beginning with the first term of this equation, for any $t > t_\lambda$,

$$\begin{aligned} & \mathbb{E}[p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q)] \\ &= \sum_{Q' \subseteq [d]} \mathbb{E}[p^*(Q') \mathbb{I}(\theta_{i_t} \geq p^*(Q') \text{ and } i_t \in Q') \mid Q = Q'] \Pr[Q = Q']. \end{aligned}$$

The event $(\theta_{i_t} \geq p^*(Q') \wedge i_t \in Q')$ is independent of the event that $Q = Q'$, so

$$\begin{aligned} \mathbb{E}[p^*(Q') \mathbb{I}(\theta_{i_t} \geq p^*(Q') \text{ and } i_t \in Q') \mid Q = Q'] &= \mathbb{E}[p^*(Q') \mathbb{I}(\theta_{i_t} \geq p^*(Q') \text{ and } i_t \in Q')] \\ &= \text{rev}(p^*(Q'), Q'). \end{aligned}$$

Therefore,

$$\begin{aligned} \sum_{t > t_\lambda} \mathbb{E}[p^*(Q) \mathbb{I}(\theta_{i_t} \geq p^*(Q) \text{ and } i_t \in Q)] &= \sum_{t > t_\lambda} \sum_{Q' \subseteq [d]} \text{rev}(p^*(Q'), Q') \Pr[Q = Q'] \\ &= \sum_{t > t_\lambda} \mathbb{E}[\text{rev}(p^*(Q), Q)]. \end{aligned} \quad (28)$$

Moving on to the second term of Equation (27), we have that for any $t > t_\lambda$,

$$\begin{aligned} & \mathbb{E}[p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q)] \\ &= \sum_{Q' \subseteq [d]} \sum_{S' \subseteq Q'} \mathbb{E}\left[\min_{i \in S'} \theta_i \cdot \mathbb{I}\left(\theta_{i_t} \geq \min_{i \in S'} \theta_i \text{ and } i_t \in Q'\right) \mid Q = Q', S_t = S'\right] \Pr[Q = Q', S_t = S']. \end{aligned} \quad (29)$$

The event that $Q = Q'$ only depends on the first t_λ timesteps and the event that $S_t = S'$ only depends on the first $t-1$ timesteps. Therefore, the event $(\theta_{i_t} \geq \min_{i \in S'} \theta_i \text{ and } i_t \in Q')$ is independent of the event $(Q = Q' \text{ and } S_t = S')$. This means that

$$\begin{aligned} & \mathbb{E}\left[\min_{i \in S'} \theta_i \cdot \mathbb{I}\left(\theta_{i_t} \geq \min_{i \in S'} \theta_i \text{ and } i_t \in Q'\right) \mid Q = Q', S_t = S'\right] \\ &= \mathbb{E}\left[\min_{i \in S'} \theta_i \cdot \mathbb{I}\left(\theta_{i_t} \geq \min_{i \in S'} \theta_i \text{ and } i_t \in Q'\right)\right] \\ &= \text{rev}\left(\min_{i \in S'} \theta_i, Q'\right). \end{aligned}$$

Combined with Equation (29), we have that

$$\begin{aligned} \mathbb{E}[p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q)] &= \sum_{Q' \subseteq [d]} \sum_{S' \subseteq Q'} \text{rev}\left(\min_{i \in S'} \theta_i, Q'\right) \Pr[Q = Q', S_t = S'] \\ &= \mathbb{E}[\text{rev}(p'_t, Q)]. \end{aligned} \quad (30)$$

Combining Equations (28) and (30), we have that

$$\mathbb{E}[Z_4] \leq \sum_{t>t_\lambda} \mathbb{E} [\text{rev}(p^*(Q), Q) - \text{rev}(p'_t, Q)]. \quad (31)$$

Next, for all $t > t_\lambda$, let \mathcal{B}_t be the event that:

1. $i_Q \in S_t$ and
2. $\text{rev}(p^*(Q), Q) - \text{rev}(p'_t, Q) \leq 4\rho_{t-1}$ (where $\rho_{t_\lambda} = 1$).

Also, let $\mathcal{C}_t = \bigcap_{s=t_\lambda+1}^t \mathcal{B}_s$. In Lemma B.2, we prove that $\Pr[\mathcal{C}_t^c] \leq \frac{1}{T}$. By Equation (31), we have that

$$\begin{aligned} \mathbb{E}[Z_4] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \text{rev}(p^*(Q), Q) - \text{rev}(p'_t, Q) \middle| \mathcal{C}_T \right] + T \Pr[\mathcal{C}_T^c] \\ &\leq \sum_{t>t_\lambda} 4\rho_{t-1} + 1 \\ &\leq 5 + 4 \sum_{t=1}^T \sqrt{\frac{\ln(dT^2)}{2t}} \\ &\leq 5 + 4\sqrt{2T \ln(dT^2)}. \end{aligned}$$

□

Claim 4.6. *If $q_{\min} \leq 2\lambda$, then $\mathbb{E}[Z_5] \leq 4\sqrt{\frac{2T}{\lambda} \ln \frac{dT^2}{\eta}} + 3$ and if $q_{\min} > 2\lambda$, $\mathbb{E}[Z_5] \leq 4\sqrt{\frac{T}{q_{\min}} \ln \frac{dT^2}{\eta}} + 2$.*

Proof. On each round $t > t_\lambda$, recall that

$$\text{LB}_{it} = \begin{cases} 0 & \text{if } \Phi_{it} = \emptyset \\ \max \left\{ \frac{1}{|\Phi_{it}|} \sum_{v \in \Phi_{it}} v - \sqrt{\frac{1}{2|\Phi_{it}|} \ln \frac{T}{\eta}}, 0 \right\} & \text{else.} \end{cases}$$

Let $j_t = \text{argmin}_{j \in S_t} \theta_j$, so $p'_t = \theta_{j_t}$. Then

$$\begin{aligned} Z_5 &= \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q) - \sum_{t>t_\lambda} p_t b_t \\ &= \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q) - \sum_{t>t_\lambda} (p'_t + (p'_t - p_t)) b_t \\ &= \sum_{t>t_\lambda} p'_t (\mathbb{I}(\theta_{i_t} \geq p'_t \text{ and } i_t \in Q) - b_t) + \sum_{t>t_\lambda} (p'_t - p_t) b_t \\ &\leq \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \wedge i_t \in Q \wedge b_t = 0) + \sum_{t>t_\lambda} (p'_t - p_t) b_t. \end{aligned}$$

Since $p_t \leq p'_t$, we have that

$$Z_5 \leq \sum_{t>t_\lambda} p'_t \mathbb{I}(\theta_{i_t} \geq p'_t \wedge i_t \in Q \wedge b_t = 0) + \sum_{t>t_\lambda} (p'_t - p_t).$$

By definition of the pricing rule, if $i_t \in S_t$, then $b_t = 1$. Therefore, if $b_t = 0$, then either $i_t \notin Q$ or $i_t \in Q \setminus S_t$. Since S_t contains every $i \in Q$ with $i > j_t$, we can conclude that if $i_t \in Q \setminus S_t$, then $\theta_{i_t} < \theta_{j_t} = p'_t$. Therefore, $\mathbb{I}(\theta_{i_t} \geq p'_t \wedge i_t \in Q \wedge b_t = 0) = 0$, which means that

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t > t_\lambda} p'_t - p_t \right].$$

Let $j'_t = \operatorname{argmin}_{j \in S_t} \text{LB}_{jt}$, which means that $p_t = \min_{i \in S_t} \{\min\{\theta_i, \text{LB}_{it}\}\} = \min\{p'_t, \text{LB}_{j'_t t}\}$. We also know that $p'_t = \theta_{j_t} \leq \theta_{j'_t}$. Therefore,

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t > t_\lambda} \max\{0, \theta_{j'_t} - \text{LB}_{j'_t t}\} \right].$$

For the remainder of our analysis, we will require the following events:

- Let \mathcal{E}_1 be the event that for all $t > t_\lambda$, $|\Phi_{i,t}| \geq \frac{1}{2}q_{\min}(t-1)$ for all $i \in S_t$. In Lemma B.4, we prove that if $q_{\min} > 2\lambda$, then $\Pr[\mathcal{E}_1^c] \leq \frac{1}{T}$.
- Similarly, let \mathcal{E}_2 be the event that for all $t > t_\lambda$, $|\Phi_{i,t}| \geq \frac{1}{4}\lambda(t-1)$ for all $i \in S_t$ such that $q_i \geq \frac{\lambda}{2}$. In Lemma B.5, we prove that if $q_{\min} \leq 2\lambda$, then $\Pr[\mathcal{E}_2^c] \leq \frac{1}{T}$.
- Let \mathcal{F} be the event that for all $i \in [d]$ such that $q_i \leq \frac{\lambda}{2}$, we have that $\bar{q}_i < \frac{3\lambda}{4}$, which means that $i \notin Q$. In Lemma B.6, we prove that $\Pr[\mathcal{F}^c] \leq \frac{1}{T}$.
- Let \mathcal{H} be the event that for all $t > t_\lambda$ and all $i \in S_t$,

$$|\Phi_{it}|\theta_i \leq \sum_{v \in \Phi_{it}} v + \sqrt{\frac{1}{2}|\Phi_{it}| \ln(dT^2)}.$$

In Lemma B.7, we prove that $\Pr[\mathcal{H}^c] \leq \frac{1}{T}$.

We now split our analysis into two cases depending on whether or not $q_{\min} > 2\lambda$. Suppose that $q_{\min} > 2\lambda$. In this case,

$$\begin{aligned} \mathbb{E}[Z_5] &\leq \mathbb{E} \left[\sum_{t > t_\lambda} \max\{0, \theta_{j'_t} - \text{LB}_{j'_t t}\} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + T \Pr[(\mathcal{E}_1 \wedge \mathcal{H})^c] \\ &\leq \mathbb{E} \left[\sum_{t > t_\lambda} \max\{0, \theta_{j'_t} - \text{LB}_{j'_t t}\} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2 \\ &= \mathbb{E} \left[\sum_{t > t_\lambda} \max \left\{ 0, \theta_{j'_t} - \frac{1}{|\Phi_{j'_t t}|} \sum_{v \in \Phi_{j'_t t}} v + \sqrt{\frac{1}{2|\Phi_{j'_t t}|} \ln \frac{1}{\eta}} \right\} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2. \end{aligned}$$

Under events \mathcal{E}_1 and \mathcal{H} ,

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t > t_\lambda} \sqrt{\frac{\ln(dT^2)}{2|\Phi_{j'_t t}|}} + \sqrt{\frac{1}{2|\Phi_{j'_t t}|} \ln \frac{1}{\eta}} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2$$

and by definition of the event \mathcal{E}_1 ,

$$\begin{aligned}\mathbb{E}[Z_5] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \sqrt{\frac{\ln(dT^2)}{2|\Phi_{j'_t}|}} + \sqrt{\frac{1}{2|\Phi_{j'_t}|} \ln \frac{1}{\eta}} \middle| \mathcal{E}_1 \wedge \mathcal{H} \right] + 2 \\ &\leq \sum_{t=2}^T \left(\sqrt{\frac{\ln(dT^2)}{q_{\min}(t-1)}} + \sqrt{\frac{1}{q_{\min}(t-1)} \ln \frac{1}{\eta}} \right) + 2 \\ &\leq 4\sqrt{\frac{T}{q_{\min}} \ln \frac{dT^2}{\eta}} + 2.\end{aligned}$$

Meanwhile, suppose that $q_{\min} < 2\lambda$. When \mathcal{F} happens, for all $t > t_\lambda$, $S_t \subseteq Q \subseteq \{i : q_i > \frac{\lambda}{2}\}$, so when both \mathcal{E}_2 and \mathcal{F} happen, $|\Phi_{i,t}| \geq \frac{1}{4}\lambda(t-1)$ for all $t > t_\lambda$ and $i \in S_t$. Therefore,

$$\begin{aligned}\mathbb{E}[Z_5] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \max \left\{ 0, \theta_{j'_t} - \text{LB}_{j'_t} \right\} \middle| \mathcal{E}_2 \wedge \mathcal{F} \wedge \mathcal{H} \right] + T \Pr[(\mathcal{E}_2 \wedge \mathcal{F} \wedge \mathcal{H})^c] \\ &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \max \left\{ 0, \theta_{j'_t} - \text{LB}_{j'_t} \right\} \middle| \mathcal{E}_2 \wedge \mathcal{F} \wedge \mathcal{H} \right] + 3 \\ &= \mathbb{E} \left[\sum_{t>t_\lambda} \max \left\{ 0, \theta_{j'_t} - \frac{1}{|\Phi_{j'_t}|} \sum_{v \in \Phi_{j'_t}} v + \sqrt{\frac{1}{2|\Phi_{j'_t}|} \ln \frac{1}{\eta}} \right\} \middle| \mathcal{E}_2 \wedge \mathcal{F} \wedge \mathcal{H} \right] + 3.\end{aligned}$$

When \mathcal{E}_2 , \mathcal{F} , and \mathcal{H} all happen,

$$\mathbb{E}[Z_5] \leq \mathbb{E} \left[\sum_{t>t_\lambda} \sqrt{\frac{\ln(dT^2)}{2|\Phi_{j'_t}|}} + \sqrt{\frac{1}{2|\Phi_{j'_t}|} \ln \frac{1}{\eta}} \middle| \mathcal{E}_2 \wedge \mathcal{F} \wedge \mathcal{H} \right] + 3$$

and by definition of $\mathcal{E}_2 \wedge \mathcal{F}$,

$$\begin{aligned}\mathbb{E}[Z_5] &\leq \mathbb{E} \left[\sum_{t>t_\lambda} \sqrt{\frac{\ln(dT^2)}{2|\Phi_{j'_t}|}} + \sqrt{\frac{1}{2|\Phi_{j'_t}|} \ln \frac{1}{\eta}} \middle| \mathcal{E}_2 \wedge \mathcal{F} \wedge \mathcal{H} \right] + 3 \\ &\leq \sum_{t=2}^T \left(\sqrt{\frac{2 \ln(dT^2)}{\lambda(t-1)}} + \sqrt{\frac{2}{\lambda(t-1)} \ln \frac{1}{\eta}} \right) + 3 \\ &\leq 4\sqrt{\frac{2T}{\lambda} \ln \frac{dT^2}{\eta}} + 3.\end{aligned}$$

□

Lemma 4.7. *For all $t > t_\lambda$, let \mathcal{A}_t be the event $\text{rev}(\theta_i, Q) \in [\check{\mu}_{i,t}, \hat{\mu}_{i,t}]$ for all $i \in S_t$. Then $\Pr[\mathcal{A}_t^c] \leq \frac{1}{T^2}$.*

Proof. Recall that $\hat{\mu}_{i,t} = \bar{\mu}_{i,t} + \rho_t$ and $\check{\mu}_{i,t} = \bar{\mu}_{i,t} - \rho_t$ with

$$\rho_t = \sqrt{\frac{\ln(dT^2)}{2(t-t_\lambda)}}.$$

We also define the related quantities for all $i \in [d]$ and all $Q' \subseteq [d]$:

$$\bar{\gamma}_{i,t}(Q') = \frac{1}{t - t_\lambda} \sum_{s=t_\lambda+1}^t \theta_i \cdot \mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q'),$$

$\hat{\gamma}_{i,t}(Q') = \bar{\gamma}_{i,t}(Q') + \rho_t$, and $\check{\gamma}_{i,t}(Q') = \bar{\gamma}_{i,t}(Q') - \rho_t$. By a Hoeffding bound, for all $Q' \subseteq [d]$ and all $i \in [d]$,

$$\Pr[\text{rev}(\theta_i, Q') \notin [\check{\gamma}_{i,t}(Q'), \hat{\gamma}_{i,t}(Q')]] \leq \frac{1}{dT^2}.$$

We claim that for any $i \in S_t$ and any $s > t_\lambda$,

$$\mathbb{I}(b_s = 1 \wedge \theta_{i_s} \geq \theta_i \wedge i_s \in Q) = \mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q), \quad (32)$$

which means that $\bar{\mu}_{i,t} = \bar{\gamma}_{i,t}(Q)$, $\hat{\mu}_{i,t} = \hat{\gamma}_{i,t}(Q)$, and $\check{\mu}_{i,t} = \check{\gamma}_{i,t}(Q)$. To see why, if $b_s = 1$, then clearly Equation (32) holds. Otherwise, suppose $b_s = 0$, in which case $\mathbb{I}(b_s = 1 \wedge \theta_{i_s} \geq \theta_i \wedge i_s \in Q) = 0$. Then $i_s \notin S_s$ because any buyer in S_s will always buy by definition of the pricing rule. Let $j_s = \min\{j \in S_s\}$. Since S_s contains every element in Q larger than j_s , we know that either:

1. $i_s \notin Q$, in which case $\mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q) = 0$, or
2. $i_s \in Q$ but $i_s \notin S_s$, which means that $\theta_{i_s} < \theta_{j_s}$. Since $i \in S_t$, it must be that $i \in S_s$, so $\theta_{i_s} < \theta_{j_s} \leq \theta_i$. In this case, $\mathbb{I}(\theta_{i_s} \geq \theta_i \wedge i_s \in Q) = 0$ as well.

Therefore, Equation (32) holds.

The fact that $\bar{\mu}_{i,t} = \bar{\gamma}_{i,t}(Q)$, $\hat{\mu}_{i,t} = \hat{\gamma}_{i,t}(Q)$, and $\check{\mu}_{i,t} = \check{\gamma}_{i,t}(Q)$ for all $i \in S_t$ implies that

$$\begin{aligned} \Pr[\mathcal{A}_t^c] &= \Pr(\exists i \in S_t \text{ s.t. } \text{rev}(\theta_i, Q) \notin [\check{\mu}_{i,t}, \hat{\mu}_{i,t}]) \\ &\leq \Pr(\exists i \in [d] \text{ s.t. } \text{rev}(\theta_i, Q) \notin [\check{\gamma}_{i,t}(Q), \hat{\gamma}_{i,t}(Q)]) \\ &\leq \sum_{i=1}^d \Pr(\text{rev}(\theta_i, Q) \notin [\check{\gamma}_{i,t}(Q), \hat{\gamma}_{i,t}(Q)]). \end{aligned} \quad (33)$$

The set Q is a random variable, so we must condition on it to bound Equation (33):

$$\begin{aligned} &\Pr(\text{rev}(\theta_i, Q) \notin [\check{\gamma}_{i,t}(Q), \hat{\gamma}_{i,t}(Q)]) \\ &= \sum_{Q' \subseteq [d]} \Pr(\text{rev}(\theta_i, Q') \notin [\check{\gamma}_{i,t}(Q'), \hat{\gamma}_{i,t}(Q')] \mid Q = Q') \Pr[Q = Q']. \end{aligned}$$

Since the event that $Q = Q'$ and the event that $\text{rev}(\theta_i, Q') \notin [\check{\gamma}_{i,t}(Q'), \hat{\gamma}_{i,t}(Q')]$ depend on disjoint timesteps, the two events are independent. Therefore,

$$\begin{aligned} \Pr(\text{rev}(\theta_i, Q) \notin [\check{\gamma}_{i,t}(Q), \hat{\gamma}_{i,t}(Q)]) &= \sum_{Q' \subseteq [d]} \Pr(\text{rev}(\theta_i, Q') \notin [\check{\gamma}_{i,t}(Q'), \hat{\gamma}_{i,t}(Q')]) \Pr[Q = Q'] \\ &\leq \frac{1}{dT^2} \sum_{Q' \subseteq [d]} \Pr[Q = Q'] \\ &= \frac{1}{dT^2}, \end{aligned}$$

so the result follows from Equation (33). \square

The next lemma shows that for more common types with $q_i \geq \lambda$, the fraction of times \bar{q}_i that that type appears during Algorithm 1 is large enough that i is added to Q .

Lemma B.1. *Let \mathcal{G} be the event that for all i such that $q_i \geq \lambda$, we have that $\bar{q}_i \geq \frac{3\lambda}{4}$. Then $\Pr[\mathcal{G}^c] \leq \frac{1}{T}$.*

Proof. Fix an index i such that $q_i \geq \lambda$. Then

$$\Pr\left[\bar{q}_i < \frac{3\lambda}{4}\right] = \Pr\left[\sum_{t=1}^{t_\lambda} \mathbb{I}(i_t = i) < \lambda t_\lambda \cdot \frac{3}{4}\right] \leq \exp\left(-\frac{\lambda t_\lambda}{32}\right) \leq \frac{1}{dT}.$$

The lemma follows by a union bound over all $i \in [d]$. \square

The next lemma proves that the expected revenue (with respect to agents in Q) of the smallest active price $\min\{\theta_i : i \in S_t\}$ is converging to the optimal revenue $\text{rev}(p^*(Q), Q)$ as t grows. Later in the analysis, we will show—at a high level—that since the algorithm sets a price within a neighborhood of $\min\{\theta_i : i \in S_t\}$, its revenue is converging to that of $p^*(Q)$. For this next lemma, recall that $p^*(Q) = \theta_{i_Q}$ for some $i_Q \in Q$. The proof is similar to that of standard successive arm elimination algorithms [e.g., Zhao and Chen, 2020].

Lemma B.2. *For all $t > t_\lambda$, let $j_t = \min\{j \in S_t\}$. Let \mathcal{B}_t be the event that:*

1. $i_Q \in S_t$ and
2. $\text{rev}(p^*(Q), Q) - \text{rev}(\theta_{j_t}, Q) \leq 4\rho_{t-1}$ (where $\rho_{t_\lambda} = 1$).

Also, let $\mathcal{C}_t = \bigcap_{s=t_\lambda+1}^t \mathcal{B}_s$. Then $\Pr[\mathcal{C}_t^c] \leq \frac{1}{T}$.

Proof. We begin by partitioning \mathcal{C}_t^c into the disjoint events

$$\mathcal{C}_t^c = \mathcal{C}_{t_\lambda+1}^c \cup (\mathcal{C}_{t_\lambda+1} \cap \mathcal{B}_{t_\lambda+2}^c) \cup \dots \cup (\mathcal{C}_{t-1} \cap \mathcal{B}_t^c).$$

Since these events are disjoint,

$$\Pr[\mathcal{C}_t^c] = \Pr[\mathcal{C}_{t_\lambda+1}^c] + \Pr[\mathcal{C}_{t_\lambda+1} \cap \mathcal{B}_{t_\lambda+2}^c] + \dots + \Pr[\mathcal{C}_{t-1} \cap \mathcal{B}_t^c]. \quad (34)$$

Beginning with the first summand, $\Pr[\mathcal{C}_{t_\lambda+1}^c] = \Pr[\mathcal{B}_{t_\lambda+1}^c] = 0$ because $S_{t_\lambda} = Q$, so $i_Q \in S_{t_\lambda}$, and $4\rho_{t_\lambda} > 1$.

Next, for $s > t_\lambda + 1$,

$$\Pr[\mathcal{C}_{s-1} \cap \mathcal{B}_s^c] = \Pr\left[\bigcap_{s'=t_\lambda+1}^{s-1} \mathcal{B}_{s'} \cap \mathcal{B}_s^c\right] \leq \Pr[\mathcal{B}_{s-1} \cap \mathcal{B}_s^c]. \quad (35)$$

We will prove that $\mathcal{B}_{s-1} \cap \mathcal{B}_s^c$ implies \mathcal{A}_{s-1}^c , which will allow us to apply Lemma 4.7.

Claim B.3. *The event $\mathcal{B}_{s-1} \cap \mathcal{B}_s^c$ implies \mathcal{A}_{s-1}^c .*

Proof. Proof of Claim B.3] First suppose \mathcal{B}_{s-1} happens and $i_Q \notin S_s$ (so \mathcal{B}_s^c happens). Since \mathcal{B}_{s-1} happens, we know that $i_Q \in S_{s-1}$ but since $i_Q \notin S_s$, it must be that i_Q was eliminated at the end of round $s-1$. This means that $\hat{\mu}_{i_Q, s-1} < \max_{k \in S_{s-1}} \check{\mu}_{k, s-1}$. Let $k' = \operatorname{argmax}_{k \in S_{s-1}} \check{\mu}_{k, s-1}$. Then

$$\begin{aligned} \operatorname{rev}(p^*(Q), Q) - \check{\mu}_{i_Q, s-1} &\geq \operatorname{rev}(\theta_{k'}, Q) - \check{\mu}_{i_Q, s-1} \\ &= \operatorname{rev}(\theta_{k'}, Q) - \hat{\mu}_{i_Q, s-1} + 2\rho_{s-1} \\ &> \operatorname{rev}(\theta_{k'}, Q) - \check{\mu}_{k', s-1} + 2\rho_{s-1}. \end{aligned} \quad (36)$$

Suppose that $\operatorname{rev}(\theta_{k'}, Q) \geq \check{\mu}_{k', s-1}$. Then Equation (36) implies that

$$2\rho_{s-1} < \operatorname{rev}(p^*(Q), Q) - \check{\mu}_{i_Q, s-1} = \operatorname{rev}(p^*(Q), Q) - (\hat{\mu}_{i_Q, s-1} - 2\rho_{s-1})$$

so $\operatorname{rev}(p^*(Q), Q) > \hat{\mu}_{i_Q, s-1}$. Therefore, either $\operatorname{rev}(\theta_{k'}, Q) < \check{\mu}_{k', s-1}$ or $\operatorname{rev}(p^*(Q), Q) > \hat{\mu}_{i_Q, s-1}$, which means that \mathcal{A}_{s-1}^c happens.

Meanwhile, suppose \mathcal{B}_{s-1} happens and $i_Q \in S_s$ but $\operatorname{rev}(p^*(Q), Q) - \operatorname{rev}(\theta_{j_s}, Q) > 4\rho_{s-1}$ (so \mathcal{B}_s^c happens). Then

$$\operatorname{rev}(p^*(Q), Q) - \check{\mu}_{i_Q, s-1} + \hat{\mu}_{j_s, s-1} - \operatorname{rev}(\theta_{j_s}, Q) > \hat{\mu}_{j_s, s-1} - \check{\mu}_{i_Q, s-1} + 4\rho_{s-1}. \quad (37)$$

Again, let $k' = \operatorname{argmax}_{k \in S_{s-1}} \check{\mu}_{k, s-1}$. Since $j_s \in S_s$, it must be that $\hat{\mu}_{j_s, s-1} \geq \check{\mu}_{k', s-1}$, or else j_s would have been eliminated at the end of round $s-1$. Combining this fact with Equation (37), we have that

$$\operatorname{rev}(p^*(Q), Q) - \check{\mu}_{i_Q, s-1} + \hat{\mu}_{j_s, s-1} - \operatorname{rev}(\theta_{j_s}, Q) > \check{\mu}_{k', s-1} - \check{\mu}_{k', s-1} + 4\rho_{s-1} = 4\rho_{s-1}.$$

This means that either:

1. $2\rho_{s-1} < \operatorname{rev}(p^*(Q), Q) - \check{\mu}_{i_Q, s-1} = \operatorname{rev}(p^*(Q), Q) - \hat{\mu}_{i_Q, s-1} + 2\rho_{s-1}$, or in other words $\hat{\mu}_{i_Q, s-1} < \operatorname{rev}(p^*(Q), Q)$, meaning \mathcal{A}_{s-1}^c happens, or
2. $2\rho_{s-1} < \hat{\mu}_{j_s, s-1} - \operatorname{rev}(\theta_{j_s}, Q) = \check{\mu}_{j_s, s-1} + 2\rho_{s-1} - \operatorname{rev}(\theta_{j_s}, Q)$, or in other words, $\operatorname{rev}(\theta_{j_s}, Q) < \check{\mu}_{j_s, s-1}$, meaning \mathcal{A}_{s-1}^c happens.

Therefore, the claim holds. \square

Claim B.3, Equation (35), and Lemma 4.7 imply that $\Pr[\mathcal{C}_{s-1} \cap \mathcal{B}_s^c] \leq \Pr[\mathcal{A}_{s-1}^c] \leq \frac{1}{T^2}$, so by Equation (34), we have that $\Pr[\mathcal{C}_t^c] < \frac{1}{T}$. \square

The next lemma will prove that for all rounds $t > t_\lambda$ of Algorithm 2 and all active types $i \in S_t$, there are a non-trivial number of reviews by buyers of type i . The following lemma holds when $q_{\min} > 2\lambda$, and Lemma B.5 holds when $q_{\min} \leq 2\lambda$.

Lemma B.4. *Suppose that $q_{\min} > 2\lambda$. Let \mathcal{E}_1 be the event that on each round $t > t_\lambda$, $|\Phi_{i,t}| \geq \frac{1}{2}q_{\min}(t-1)$ and all $i \in S_t$. Then $\Pr[\mathcal{E}_1^c] \leq \frac{1}{T}$.*

Proof. Fix any $t > t_\lambda$. We will show that

$$\Pr \left[\exists i \in S_t \text{ such that } |\Phi_{i,t}| < \frac{1}{2}q_{\min}(t-1) \right] \leq \frac{1}{T^2}.$$

By definition, $|\Phi_{i,t}| = \sum_{s=1}^{t-1} \mathbb{I}(b_s = 1 \wedge i_s = i)$. If $|\Phi_{i,t}|$ were equal to $\sum_{s=1}^{t-1} \mathbb{I}(i_s = i)$, then the claim would hold immediately by a Chernoff bound. However, we do not know at each round s whether $i_s = i$ provided the buyer does not make a purchase. Therefore, we also define the random variable $X_{i,t} = \sum_{s=1}^{t-1} \mathbb{I}(i_s = i)$. We claim that for all $i \in S_t$, $|\Phi_{i,t}| = X_{i,t}$. This is because we know that $i \in S_s$ for all $s \leq t$ and by definition of the pricing rule, if $i_s = i$, then $b_s = 1$.

Therefore,

$$\begin{aligned}
& \Pr \left[\exists t > t_\lambda, \exists i \in S_t \text{ such that } |\Phi_{i,t}| < \frac{1}{2} q_{\min}(t-1) \right] \\
&= \Pr \left[\exists t > t_\lambda, \exists i \in S_t \text{ such that } X_{i,t} < \frac{1}{2} q_{\min}(t-1) \right] \\
&\leq \Pr \left[\exists t > t_\lambda, \exists i \in [d] \text{ such that } X_{i,t} < \frac{1}{2} q_{\min}(t-1) \right] \\
&\leq \sum_{i=1}^d \sum_{t=t_\lambda+1}^T \Pr \left[X_{i,t} < \frac{1}{2} q_{\min}(t-1) \right]. \tag{38}
\end{aligned}$$

By a Chernoff bound,

$$\begin{aligned}
\Pr \left[X_{i,t} \leq \frac{1}{2} q_{\min}(t-1) \right] &\leq \Pr \left[X_{i,t} \leq \frac{1}{2} q_i(t-1) \right] \\
&\leq \exp \left(-\frac{q_i(t-1)}{8} \right) \\
&\leq \exp \left(-\frac{q_{\min}(t-1)}{8} \right) \\
&\leq \exp \left(-\frac{\lambda(t-1)}{4} \right) \\
&\leq \frac{1}{dT^2}.
\end{aligned}$$

The lemma now follows from Equation (38). \square

We now prove a similar result for the case where $q_{\min} \leq 2\lambda$.

Lemma B.5. *Suppose that $q_{\min} \leq 2\lambda$. Let \mathcal{E}_2 be the event that for all $t > t_\lambda$, $|\Phi_{i,t}| \geq \frac{1}{4}\lambda(t-1)$ for all $i \in S_t$ such that $q_i \geq \frac{\lambda}{2}$. Then $\Pr[\mathcal{E}_2^c] \leq \frac{1}{T}$.*

Proof. Let $Q_0 = \{i : q_i \geq \frac{\lambda}{2}\}$. Fix any $t > t_\lambda$. We will show that

$$\Pr \left[\exists i \in S_t \cap Q_0 \text{ such that } |\Phi_{i,t}| < \frac{1}{4}\lambda(t-1) \right] \leq \frac{1}{T^2}.$$

By definition, $|\Phi_{i,t}| = \sum_{s=1}^{t-1} \mathbb{I}(b_s = 1 \wedge i_s = i)$. As in the proof of Lemma B.4, we define the random variable $X_{i,t} = \sum_{s=1}^{t-1} \mathbb{I}(x_s = e_i)$. As in that proof, for all $i \in S_t$, $|\Phi_{i,t}| = X_{i,t}$ (this is because we know that $i \in S_s$ for all $s \leq t$ and by definition of the pricing rule, if $i_s = i$, then $b_s = 1$).

Therefore,

$$\begin{aligned}
& \Pr \left[\exists t > t_\lambda, \exists i \in S_t \cap Q_0 \text{ such that } |\Phi_{i,t}| < \frac{1}{4}\lambda(t-1) \right] \\
&= \Pr \left[\exists t > t_\lambda, \exists i \in S_t \cap Q_0 \text{ such that } X_{i,t} < \frac{1}{4}\lambda(t-1) \right] \\
&\leq \Pr \left[\exists t > t_\lambda, \exists i \in Q_0 \text{ such that } X_{i,t} < \frac{1}{4}\lambda(t-1) \right] \\
&\leq \sum_{i \in Q_0} \sum_{t=t_{\lambda+1}}^T \Pr \left[X_{i,t} < \frac{1}{4}\lambda(t-1) \right]. \tag{39}
\end{aligned}$$

By a Chernoff bound, for any $i \in Q_0$,

$$\begin{aligned}
\Pr \left[X_{i,t} < \frac{1}{4}\lambda(t-1) \right] &\leq \Pr \left[X_{i,t} < \frac{1}{2}q_i(t-1) \right] \\
&\leq \exp \left(-\frac{q_i(t-1)}{8} \right) \\
&\leq \exp \left(-\frac{\lambda(t-1)}{16} \right) \\
&\leq \frac{1}{dT^2}.
\end{aligned}$$

The lemma therefore follows from Equation (39). \square

We next observe that for all very rare types with $q_i \leq \lambda/2$, the fraction of times \bar{q}_i that that type appears during Algorithm 1 is small. Therefore, i is not added to the set Q and is ignored for the remainder of the algorithm.

Lemma B.6. *Let \mathcal{F} be the event that for all $i \in [d]$ such that $q_i \leq \frac{\lambda}{2}$, we have that $\bar{q}_i \leq \frac{3\lambda}{4}$. Then $\Pr[\mathcal{F}^c] \leq \frac{1}{T}$.*

Proof. Fix an index i such that $q_i \leq \frac{\lambda}{2}$. Then

$$\Pr \left[\bar{q}_i \geq \frac{3\lambda}{4} \right] = \Pr \left[\sum_{t=1}^{t_\lambda} \mathbb{I}(i_t = i) \geq \frac{\lambda t_\lambda}{2} \cdot \frac{3}{2} \right] \leq \exp \left(-\frac{\lambda t_\lambda}{24} \right) = \frac{1}{dT}.$$

The lemma the follows by a union bound over all $i \in [d]$. \square

Our final lemma proves that for all active types $i \in S_t$, the average reviews of agents with this type is close to the true *ex-ante* value θ_i . This helps us ensure that the price we set is not too low.

Lemma B.7. *Let \mathcal{H} be the event that for all $t > t_\lambda$ and all $i \in S_t$,*

$$|\Phi_{it}| \theta_i \leq \sum_{v \in \Phi_{it}} v + \sqrt{\frac{1}{2} |\Phi_{it}| \ln(dT^2)}.$$

Then $\Pr[\mathcal{H}^c] \leq \frac{1}{T}$.

Proof. Fix any $t > t_\lambda$. Let v_1, \dots, v_{t-1} be the buyers' *ex-post* values (which are defined even if the buyer didn't buy on a particular round s as $v_s \sim \mathcal{D}_{i_s}$). For each $i \in [d]$, let $R_{it} = \{s < t : i_s = i\}$ be the set of rounds in which the buyer had type i . Since any buyer $i \in S_s$ will buy if $i_s = i$, we have that $|\Phi_{it}| = |R_{it}|$ and

$$\sum_{v \in \Phi_{it}} v = \sum_{s \in R_{it}} v_s.$$

Therefore,

$$\begin{aligned} & \Pr \left[\exists i \in S_t \text{ such that } |\Phi_{it}| \theta_i > \sum_{v \in \Phi_{it}} v + \sqrt{\frac{1}{2} |\Phi_{it}| \ln(dT^2)} \right] \\ & \leq \Pr \left[\exists i \in [d] \text{ such that } |R_{it}| \theta_i > \sum_{s \in R_{it}} v_s + \sqrt{\frac{1}{2} |R_{it}| \ln(dT^2)} \right] \\ & \leq \sum_{i=1}^d \Pr \left[|R_{it}| \theta_i > \sum_{s \in R_{it}} v_s + \sqrt{\frac{1}{2} |R_{it}| \ln(dT^2)} \right] \\ & = \sum_{i=1}^d \sum_{R \subseteq [t-1]} \Pr \left[|R| \theta_i > \sum_{s \in R} v_s + \sqrt{\frac{1}{2} |R| \ln(dT^2)} \mid R_{it} = R \right] \Pr[R_{it} = R]. \end{aligned} \quad (40)$$

For any $s \in R$, $\mathbb{E}[v_s \mid R_{it} = R] = \theta_i$. Therefore,

$$\Pr \left[|R| \theta_i > \sum_{s \in R} v_s + \sqrt{\frac{1}{2} |R| \ln(dT^2)} \mid R_{it} = R \right] \leq \frac{1}{dT^2}.$$

The lemma therefore follows from Equation (40) and a union bound over all rounds $t > t_\lambda$. \square