Learning to Mitigate AI Collusion on Economic Platforms

Anonymous Author(s) Affiliation Address email

Abstract

1	Algorithmic pricing on online e-commerce platforms raises the concern of tacit
2	collusion, where reinforcement learning algorithms learn to set collusive prices
3	in a decentralized manner and through nothing more than profit feedback. This
4	raises the question as to whether collusive pricing can be prevented through the
5	design of suitable "buy boxes," i.e., through the design of the rules that govern
6	the elements of e-commerce sites that promote particular products and prices to
7	consumers. In this paper, we demonstrate that reinforcement learning (RL) can also
8	be used by platforms to learn buy box rules that are effective in preventing collusion
9	by RL sellers. For this, we adopt the methodology of <i>Stackelberg POMDPs</i> , and
10	demonstrate success in learning robust rules that continue to provide high consumer
11	welfare together with sellers employing different behavior models or having out-of-
12	distribution costs for goods.

13 **1 Introduction**

The last decade has witnessed a dramatic shift of trading from retailers to online e-commerce platforms 14 such as Amazon and Alibaba. In these platforms, sellers are increasingly using algorithms to set prices. 15 Algorithmic pricing can be beneficial for market efficiency, enabling sellers to quickly react to market 16 changes and also in enabling price competition. At the same time, the U.S. Federal Trade Commission 17 (FTC) U.S. Federal Trade Commission (2018) and European Commission (The Organisation for 18 Economic Co-operation and Development, 2017) have raised concerns that algorithmic pricing may 19 facilitate collusive behaviors. Calvano et al. (2020a) support these concerns through a study of pricing 20 agents in a simulated platform economy, and show that commonly used reinforcement-learning (RL) 21 algorithms learn to initiate and sustain collusive behaviors. Assad et al. (2020) also provide empirical 22 support for algorithmic collusion in a study of Germany's retail gas stations, showing an association 23 between algorithmic pricing and an increase in price markups. As highlighted by Calvano et al. 24 (2020b), these kinds of collusive behaviors are unlikely to be a violation of antitrust law, as they are 25 learned responses to profit signals and not the result of explicit agreements. 26

One can try to prevent algorithmic collusion by introducing suitable rules by which platforms can choose which sellers to promote to buyers, thus promoting competition. Could Amazon's "buy box algorithm," for example, play this role in the future, in determining for a given consumer search which products and prices to highlight to a consumer? Responding to this, Johnson et al. (2020) design hand-crafted rules that succeed in hindering collusion between RL algorithms. At the same time, their rules introduce the undesirable effect of limiting consumers to a single seller, and there remains potential for more effective interventions.

In this paper, we demonstrate for the first time how RL can also be used defensively by a platform
 to automatically design rules that promote consumer welfare and prevent collusive pricing. This is
 a problem of multi-agent learning, with the interaction between the platform and sellers modeled

as a *Stackelberg game* (Fudenberg and Tirole, 1991). The leader is the platform designer and sets

the platform rules and the sellers respond, using RL to set prices given these rules. We introduce the

class of *threshold platform rules*, and formally show that this class contains rules that approximately maximize consumer surplus in a unique subgame perfect equilbrium. At the same time, this class of

threshold rules is fragile to unexpected deviations by sellers, for example caused by cost perturbations.

The role of RL on the part of the platform is to learn rules with similar performance that are also

43 more robust.

To solve the Stackelberg problem, we make use of the *Stackelberg partially observable Markov decision process* (POMDP) framework (Brero et al.) 2021a), which defines an episode structure of a POMDP such that the RL algorithm representing the leader will learn to optimize its reward (here, consumer surplus) given that its rules cause re-equilibration on the part of the followers (here, the sellers who use Q-learning algorithms to set prices). The Stackelberg POMDP framework is well formed as long as the re-equilibration behavior of sellers can be modeled through Markovian dynamics, as is the case with Q-learning.

We show successful results in learning effective platform policies that outperform handcrafted 51 rules (Johnson et al.) 2020). This demonstrates how the Stackelberg POMDP framework can be 52 successfully applied in settings where followers play repeated games, and their strategies are also 53 policies trained via reinforcement learning algorithms. We then show how our threshold platform 54 rules allow us to obtain a similar learning performance when training the platform policy "in the 55 wild," i.e., without accessing the sellers' private information. With this, we demonstrate how the 56 Stackelberg POMDP framework can be applied in more general learning scenarios than the offline 57 learning ones for which it was originally designed. Finally, we show how the platform rules learned 58 via our Stackelberg POMDP framework continue to be effective when market conditions change, for 59 example as the result of a change to the cost structure of sellers. 60

Further related work. Zheng et al. (2022); Tang (2017); Shen et al. (2020); Brero et al. (2021b) 61 make use of RL to optimize different economic systems (including matching markets, internet 62 advertising, tax policies, and auctions) under strategic agents' responses. Unlike our work, these 63 methods do not leverage the designer's commitment power or the Stackelberg structure of the induced 64 game. Brero et al. (2021a) introduce and study the Stackelberg POMDP framework for a very different 65 setting than that of the present paper: the design of sequential price auctions,¹ Abada and Lambin 66 (2020) study collusion by RL pricing in markets for electric power, and use machine learning by a 67 regulator agent for the mitigation of collusion, albeit without a Stackelberg framing (and without 68 success, leading to lower welfare than the collusive outcome). The broader research program on 69 differentiable economics uses representation learning for optimal economic design (Duetting et al., 70 2019; Shen et al., 2019; Kuo et al., 2020; Tacchetti et al., 2019; Rahme et al., 2021a; Curry et al., 71 2022; Rahme et al., 2021b; Curry et al., 2020; Peri et al., 2021; Weissteiner and Seuken, 2020); this 72 work avoids the need for Stackelberg design by emphasizing the use of direct, incentive-compatible 73 mechanisms. Also related is *empirical mechanism design* (Areyan Viqueira et al., 2019; Vorobeychik 74 et al., 2006; Brinkman and Wellman, 2017), which applies empirical game theory to search for the 75 equilibria of mechanisms with a set of candidate strategies (Wellman, 2006; Kiekintveld and Wellman, 76 2008; Jordan et al., 2010); see also Bünz et al. (2018) for the design of iterative auctions. 77

78 2 Preliminaries

79 Seller Competition Model. There is a set of sellers $\mathcal{N} = \{1, \ldots, n\}$, each of whom sells a 80 differentiated product on an economic platform. Each seller has the same (private) marginal cost 81 c > 0 for producing one unit of its product. Sellers interact with each other repeatedly over time in 82 setting prices and selling goods. At each time period, $t = 0, 1, \ldots$, each seller *i* observes all past 83 prices, and sets a price $p_{i,t} \ge 0$. We let $p_t = (p_{1,t}, \ldots, p_{n,t})$ denote a generic price profile quoted at

¹The only other method we know for Stackelberg learning in stochastic games with provably guarantees solves for a single follower (Mishra et al.) (2020); see also Mguni et al. (2019); Cheng et al. (2017); Shi et al. (2020) and Shu and Tian (2019), and Tharakunnel and Bhattacharyya (2007) for a partial convergence result for a static game with two followers. For other convergence results for single-follower, static, and often zero-sum games see Li et al. (2019); Sengupta and Kambhampati (2020); Xu et al. (2021); Fiez et al. (2020); Jin et al. (2020). For multi-follower static games, Wang et al. (2022) make use of a differentiable relaxation of follower best-response behavior together with a subroutine to solve an optimization problem for follower behavior.

time t. The platform sets the rules of a buy box that govern, in each period t, which set $N_t \subseteq N$ of

sellers are available. Consumers can only buy from these sellers and others forfeit sales. There is
 also an outside option, indexed by 0, which provides each consumer with a fallback choice with zero
 utility.

Following Johnson et al. (2020), competition between sellers for consumer demand is modeled 88 through the standard *logit model* of consumer choice. For this, seller i has quality index $\alpha_i > 0$, 89 this providing horizontal differentiation across products, and the outside good has quality index 90 $\alpha_0 > 0$. In the logit model, each consumer samples $\zeta_0, \zeta_1, ..., \zeta_n$, independently from a type I extreme 91 value distribution with scale parameter $\mu > 0$, for each product and the outside option, with utility 92 $\alpha_i + \zeta_i - p_{i,t}$ for product i, and $\alpha_0 + \zeta_0$ for the outside option. Considering a unit mass of consumers in 93 period t, seller $i \in \mathcal{N}_t$ receives fractional demand $D_i(p_t; \mathcal{N}_t) = \exp((\alpha_i - p_{i,t})/\mu)/\lambda(p_t; \mathcal{N}_t)$, where $\lambda(p_t; \mathcal{N}_t) = \sum_{j \in \mathcal{N}_t} \exp((\alpha_j - p_{j,t})/\mu) + \exp(\alpha_0/\mu)$, and any seller $i \notin \mathcal{N}_t$ has zero demand. Scale 94 95 parameter $\mu > 0$ serves to control the extent of horizontal differentiation, with no differentiation and 96 perfect substitutes obtained as $\mu \to 0$. The total consumer surplus is $U(p_t; \mathcal{N}_t) = \mu \cdot \log[\lambda(p_t; \mathcal{N}_t)]$, 97 and is maximized with minimum prices and all sellers displayed (so consumers have a full choice of 98 products). Seller i's profit ρ_i in period t is $\rho_i(p_t; \mathcal{N}_t) = (p_{i,t} - c) \cdot D_i(p_t; \mathcal{N}_t)$, and its per-unit profit 99 multiplied by demand. 100

Reinforcement learning by sellers. In a single-agent Markov decision process (MDP), an agent 101 faces a sequential decision problem under uncertainty. At each step t, the agent observes a state 102 variable $s_t \in S$ and chooses an action $a_t \in A$. Upon action a_t in state s_t , the agent obtains 103 reward $r(s_t, a_t)$, and the environment moves to state s_{t+1} according to $p(s_{t+1}|s_t, a_t)$. We let 104 $\tau = (s_0, a_0, ..., s_T, a_T)$ denote a state-action trajectory determined by executing policy policy 105 $\pi : S \to A$, and $p_{\pi}(\tau)$ denote the probability of trajectory τ . The optimal policy π^* solves 106 $\pi^* \in \operatorname{argmax}_{\pi} E_{\tau \sim p_{\pi}(\tau)} [\sum_{t=0}^{T} \delta^t r(s_t, a_t)]$, where $\delta \in [0, 1]$ is the discount factor and time-horizon T can be finite or infinite. In a *partially-observable MDP* (POMDP), the policy π cannot access 107 108 state s_t but only observation o_t sampled from $p(o_t|s_t)$. A multi-agent MDP (Boutilier, 1996) for 109 n agents has states S common to all agents and a set of actions A_i for each agent i. When each 110 agent i picks action $a_{i,t}$ in state s_t , the environment moves to state s_{t+1} according to a distribution 111 $p(s_{t+1}|s_t, a_{1,t}, ..., a_{n,t})$ and agent i obtains a reward $r_i(s_t, a_t)$ that depends on the joint action. We 112 follow Calvano et al. (2020a) and Johnson et al. (2020) and adopt decentralized Q-learning by sellers 113 as a positive theory for sellers in regard to their behavior in setting prices on an e-commerce platform 114 (see Appendix A). Although Q-learners may not converge, we also confirm these earlier studies in 115 showing convergence in our simulations (defined over a particular time horizon as detailed by Johnson 116 et al. (2020)). 117

118 3 The Platform Stackelberg Problem

To formalize the problem facing the platform in mitigating collusive behavior by sellers, we model 119 the interaction between the platform, which sets the rules of the buy box, and the sellers as a 120 Stackelberg game. The platform designer is the leader, and fixes the platform rules. The sellers are 121 the followers, and play an infinitely repeated game according to these rules. As discussed above, and 122 following Calvano et al. (2020a) and Johnson et al. (2020), we model the sellers' behavior through 123 decentralized Q-learning. As a result, the problem facing the platform is a *behavioral Stackelberg* 124 problem, in that the followers are modeled as Q-learners (and need not, necessarily, be playing an 125 equilibrium of the induced game). 126

The sellers. In this model, we fix the states that comprise the MDP of a seller to include the prices set by all sellers in the last period, i.e., $s_t = p_{t-1}$. We initialize s_0 to be a randomly selected price profile. The action of a seller is modeled as one of m equally-spaced points in the interval ranging from just below the sellers' cost c to just above the monopoly price p^m . At each step $t \ge 0$, each seller i selects a price $p_{i,t}$ and is rewarded by its per-period profit $\rho_i(p_t; \mathcal{N}_t)$, which depends on $p_t = (p_{1,t}, \ldots, p_{n,t})$ and the choice of which sellers \mathcal{N}_t are displayed by the platform.

The platform. To formalize the platform's problem, let $\sigma^* = (\sigma_1^*, ..., \sigma_n^*)$ denote a strategy profile selected by Q-learning on the part of sellers, in response to the platform rule, and in the long run, after a suitably large number of steps. We leave implicit here the dependence of seller strategy profile on

the platform's policy. The platform must decide in each period which sellers to display to consumers. 136 For this, we denote the platform rule as *policy* π , and we adopt for the state of the platform policy 137 the prices quoted by sellers in step t, p_t , so that the platform's policy uses these prices to select a 138 set of agents to display, with \mathcal{N}_t selected according to $\pi(p_t)$. Let $p_t^* = \sigma^*(s_t)$ denote a price profile 139 chosen under seller strategies σ^* , i.e., in response to the platform rules, and at some large enough 140 time step t^* , and let $\tau^* = (p_{t^*}^*, p_{t^*+1}^*, ..)$ denote a trajectory of prices forward from t^* . We denote 141 the distribution of these trajectories as $p_{\pi}(\tau^*)$. As above, the dependence on the platform's policy is 142 left implicit in this notation. 143

The Stackelberg problem facing the platform is to find a platform policy π that maximizes consumer surplus given the effect of this policy on the induced strategy profile of sellers.

146 **Definition 1 (Behavioral Stackelberg Problem)** The optimal platform policy solves $\pi^* \in$ 147 $\operatorname{argmax}_{\pi} CS(\pi)$, where $CS(\pi)$ is the expected sum consumer surplus when sellers follow strategy σ^* 148 forward from period t^* , i.e.,

$$CS(\pi) = \mathbb{E}_{\tau^* \sim p_{\pi}(\tau^*)} \left[\sum_{t=t^*}^{T^*} U(p_t^*; \pi(p_t^*)) \right],$$
(1)

where T^* is suitably chosen horizon and $p_{\pi}(\tau^*)$ denotes the distribution over *Q*-learning induced, seller pricing trajectories, in response to platform policy π .

151 4 Learning Optimal Platform Rules

In this section, we solve the platform's problem, in responding to Q-learning sellers, through the *Stackelberg POMDP* framework (Brero et al., 2021a). This creates a suitably defined POMDP in which the optimal policy solves the behavioral Stackelberg problem (Definition 1).

Definition 2 (Stackelberg POMDP for platform rules) *The* Stackelberg POMDP for platform rules *is a finite-horizon POMDP, where each episode has the following two phases:*

• An equilibrium phase, consisting of $n_e \ge 1$ steps. In this phase, each state s_t includes the step counter t, the sellers' current Q-matrices, and the prices p_t quoted by the agents. Observations consists of the prices quoted by the sellers ($o_t = p_t$) and policy actions determine the set of agents displayed (in their more general version, $a_t = N_t$). State transitions are determined by Q-learning, where each seller i updates its Q-matrix after being rewarded by $\rho_i(p_t; N_t)$. The policy has zero reward in every time step ($r(s_t, a_t) = 0$, for $t \le n_e$).

• A reward phase, consisting of $n_r \ge 1$ steps, each with the same actions, states, and observations as the equilibrium steps. The reward phase differs in two ways. First, the Q-matrices of sellers are not updated, and second, the platform policy now receives a non-zero reward, and this is set in each step to be equal to the consumer surplusin that step $(r(s_t, a_t) = U(p_t; \mathcal{N}_t), \text{ for } t > n_e)$.

This Stackelberg POMDP formulation is an adaptation of that provided by Brero et al. (2021a), who used it to learn sequential price mechanisms (SPMs) in the presence of communication from bidders. Here, our stage games replace SPMs, and the followers respond through Q-learning dynamics rather than no-regret algorithms. Following Brero et al. (2021a), we show the Stackelberg POMDP formulation is well-founded by showing that an optimal policy will also solve the Behavioral Stackelberg design problem of Definition I. Specifically, when the number of reward steps n_r is large enough and when $n_e \ge t^*$, the optimal policy, denoted $\pi^*_{n_e,n_r}$, for the Stackelberg POMDP with n_e equilibrium and n_r reward steps maximizes the objective in Equation (I).

Proposition 1 The optimal Stackelberg POMDP policy π_{n_e,n_r}^* , for an equilibrium phase with $n_e \ge 1$ steps and a reward phase with $n_r \ge 1$ steps, maximizes $CS(\pi)$, for seller behavior induced after n_e steps when $n_r = T^*$.

The proposition follows from the construction of the Stackelberg POMDP, especially the fact the our policy is only rewarded under the response behavior reached after n_e steps, in line with the definition of $CS(\pi)$ (see Appendix B for the full proof argument).

181 Brero et al. (2021a) use the Stackelberg POMDP framework in an "offline" environment, i.e., in 182 a simulation that assumes access, at design time, to followers' internal information. This allows them to solve their POMDP using the paradigm of *centralized training and decentralized execution* (Lowe et al., 2017). The leader policy is trained via an actor-critic deep RL algorithm, and the critic network (which estimates the sum of rewards until the end of the episode) accesses the full state during training. Only the actor network, which represents the policy, is restricted to the partial-state information.

Here, we also study the use of the Stackelberg POMDP framework to train useful leader policies "in
the wild," where the learning algorithm of the platform can only access the kind of information that an
economic platform would have based on observations of sellers. As we will empirically demonstrate,
we can successfully relax the offline learning requirements—i.e., we operate without (1) access to
sellers' private information in regard to Q-matrices and exploration rate, and (2) requiring that the
Q-matrices of sellers become frozen for the reward phase of the Stackelberg POMDP—without
affecting learning performance²

Threshold platform rules. In our experiments, we consider the class of *threshold platform rules*. These rules use the current prices set by sellers to set a price threshold above which a seller will not be displayed, with the same threshold set for all sellers.

Definition 3 (Threshold Platform Rule) A threshold platform rule sets a threshold $\tau(p_t) \ge 0$, for each price profile p_t , such that $\mathcal{N}_t = \{i \in \{1, ..., n\} : p_{i,t} \le \tau(p_t)\}$, i.e., any seller whose price is no greater than the threshold is displayed to consumers.

This class of threshold rules has a corresponding optimality result: there is a threshold rule that makes the market competitive, with all sellers displayed and consumer surplus maximized in the unique subgame perfect Nash equilibrium (SPE) of the induced continuous pricing game. Even though the pricing behaviors that arise from Q-learning need not converge to SPEs, we find empirical evidence, consistent with Johnson et al. (2020), that the seller learning dynamics invariably converge to this equilibrium. As such, this provides useful theoretical support for adopting the family of threshold platform rules by the platform learner. We have the following result:

Proposition 2 For any $\epsilon > 0$, there exists a threshold platform rule π such that $CS(\pi) \ge CS(\pi^*) - \sum_t \delta^t \epsilon$ under the unique subgame perfect Nash equilibrium (SPE) of any finitely-repeated continuous pricing game induced by platform rule π .

This proposition follows from a platform rule with a limiting threshold that is arbitrarily close to the sellers' cost c (see Appendix C for the proof). Under this rule, sellers are displayed only if their price is minimal. At the same time, this particular threshold rule is fragile, and would lead to market failure if seller costs vary in unexpected ways. By letting the threshold τ also vary with the price profile p_t , as is allowed by the family of threshold platform rules, we seek to learn milder interventions that still mitigate collusion but remain robust to variations in the costs faced by sellers in the marketplace.

217 **5 Experimental Results**

218 In this section, we evaluate our learning approach via three main experiments. We first consider performance in terms of consumer surplus, benchmarking our RL interventions against the ones 219 introduced by Johnson et al. (2020). We demonstrate the ability to learn optimal leader strategies in 220 the Stackelberg game with the followers across all the seeds we tested, significantly outperforming 221 existing interventions. We then train platform rules without access to the sellers' private information 222 ("in the wild,") and show that this is not crucial for our learning performances. We conclude by 223 testing the robustness of our interventions, adding price perturbation during training and evaluating 224 the effect on the robustness of our learned platform rules in environments where sellers have different 225 costs from those assumed during training. 226

Experimental set-up. As in Calvano et al. (2020a) and Johnson et al. (2020), we consider settings with two pricing agents with cost c = 1, quality indexes $\alpha_1 = \alpha_2 = 2$, and $\alpha_0 = 0$, and we set parameter $\mu = 0.25$ to control horizontal differentiation. The seller Q-learning algorithms are also

²We notice this is also in line with the recent findings in Fujimoto et al.] (2022) highlighting how the Bellman error minimization (for which we require environments to be Markovian) may not be a good proxy of the accuracy of the value function.

trained using discount factor $\delta = 0.95$, exploration rate $\varepsilon_t = e^{-\beta t}$ with $\beta = 1e - 5$, and learning rate $\alpha = 0.15$. We also include results for variations of this default setting in the Appendix.

We adopt five prices choices for the action of each seller, these prices ranging from just below the 232 sellers' cost to the monopoly price. Earlier work provided sellers with a choice of fifteen different 233 prices (over a similar range). We need a smaller grid in order to satisfy our computational constraints; 234 earlier work studied the effect of different, hand-designed platform rules, and did not also use RL for 235 the automated design of suitable platform rules. We also follow the choices of earlier work, and study 236 an economy with two sellers (again, for reasons of computational resources). This coarsened price 237 grid allows us to train a platform policy through Stackelberg POMDP for 50 million steps in 18 hours 238 using a single core on a Intel(R) Xeon(R) Platinum 8268 CPU @ 2.90GHz machine. 239

Learning algorithm. To train the platform policy, we start from the A2C algorithm provided by 240 Stable Baselines3 (Raffin et al., 2021, MIT License). Given that our policy is only rewarded at 241 the end of a Stackelberg POMDP episode, we configure A2C so that neural network parameters 242 are only updated after this reward phase. In this way, we guarantee that policies inducing desired 243 followers' equilibria are properly rewarded. Furthermore, to reduce variance in sellers' responses 244 due to non-deterministic policy behavior, we maintain an observation-action map throughout each 245 episode. When a new observation is encountered during the episode, the policy chooses an action 246 following the default training behavior and stores this new observation-action pair in the map. We 247 will show the importance of this variation via an ablation study that is presented in the Appendix. 248 Sellers restart the Q-learning process by re-initializing exploration rates every time the platform rules 249 change (i.e., at the beginning of every Stackelberg POMDP episode). We also show how the training 250 approach is robust to different sellers' behavior models in the Appendix, where the sellers restart the 251 learning rate asynchronously, and not necessarily at the beginning of episodes. 252

253 5.1 Platform Learning Performance

In this section, we evaluate the performance of our learned platform policies. For this, we train our 254 policies for 50 million steps in total. We set up the Stackelberg POMDP environment using 50k 255 equilibrium steps and 30 reward steps.³ In these initial experiments, we train our policies using the 256 centralized training-decentralized execution paradigm as used for this Stackelberg learning problem 257 by Brero et al. (2021a), giving the critic network access to the sellers' learning information (i.e., 258 Q-tables and exploration rates). We relax this below in studying the robustness of the computational 259 framework to online training ("in the wild.") We consider the following interventions on behalf of the 260 platform designer: 261

• *No intervention:* Sellers are always displayed, no matter the price they quote. To derive this baseline, we run our Q learning dynamics until convergence (as described in Johnson et al. (2020)) for each seed and then average the surplus at final strategies.

• *PDP*: We test *price-directed prominence*, a platform intervention introduced by Johnson et al. (2020). Here, the platform only displays the seller who quotes the lower price (breaking ties at random), thus enhancing competition. As for *no intervention*, we compute the performance of *PDP* by averaging consumer surplus after Q-learning dynamics converge.

• *DPDP: Dynamic price-directed prominence* is another intervention introduced by Johnson et al. (2020), which also conditions the choice of the (unique) displayed seller on past prices. Under this intervention, quoting prices equal to cost is a subgame perfect equilibrium of the induced game (under suitable discount factors). As for the previous baselines, we compute the performance of *DPDP* by averaging consumer surplus after Q-learning dynamics converge.

• No State RL: Here we use the Stackelberg POMDP methodology to train a platform policy that does not use prices p_t to determine the threshold at which to admit each seller to the buy box (thus, "no state"),⁴ Here, Q-learning is restarted whenever a Stackelberg POMDP episode begins.

•*No Stackelberg No State RL:* A variation on "No State RL" that does not use the Stackelberg MDP methodology. Rather, the platform and sellers each follow decentralized learning, and the platform

³See the Appendix for a discussion around parameter selection.

⁴This class of policies already includes the optimal policy described in the proof of Proposition 2



Figure 1: Learning performance of No State RL and State-based RL compared with different baselines. The results are averaged over 10 runs and shaded regions show 95% conf. intervals. The No-Stackelberg interventions are displayed on the left, the Stackelberg ones are on the right.

279 receives a consumer surplus reward at every step. Q-learning is restarted after the same number of280 steps that are used in a Stackelberg MDP episode.

• *State-based RL:* Here we use the Stackelberg POMDP methodology to train a platform policy that sets a threshold at which to admit each seller as a function of the price profile quoted by the sellers (thus, "state-based"). This is the full class of threshold platform rules. Here, Q-learning is restarted whenever a Stackelberg POMDP episodes begins.

No Stackelberg State-based RL: A variation on "State-based RL" that does not use the Stackelberg
 MDP methodology. Rather, the platform and sellers each follow decentralized learning, and the
 platform receives a consumer surplus reward at every step. Q-learning is restarted after the same
 number of steps that are used in a Stackelberg MDP episode.

Figure I shows the consumer surplus that is realized under these different interventions. First, we 289 confirm the results of Johnson et al. (2020), and see consumer surplus improvements from both 290 *PDP* and *DPDP* compared to *No intervention*, with *DPDP* outperforming *PDP*. At the same time, 291 292 the no Stackelberg baselines are not able to outperform DPDP, confirming the benefits of using learning methodologies that exploit the leader-follower structure of our game. Indeed, our RL 293 interventions based on the Stackelberg framework dramatically improve consumer surplus, driving it 294 to (approximately, in the state-based scenario) its maximal level. In our setting, this optimal level for 295 surplus is approximately 0.94. This is confirmed by the fact that, for both No State and State-based 296 RL, all sellers are displayed and they invariably quote minimum prices at the end of training. This is 297 the optimal (i.e., surplus maximizing) seller behavior, confirming the effectiveness of the Stackelberg-298 based learning methodology in finding an optimal leader strategy given the Q-learning behavior of 299 sellers. It is easier for *No State RL* to reach the optimal performance since its class of policies is 300 much smaller than the class considered by *State-based RL*. However, as we will see in Section 5.3 301 the state-based policy is more flexible and is robust to the case that the cost basis changes for sellers 302 while No State RL is not. 303

304 5.2 Learning in the Wild

We now test the performance of the Stackelberg POMDP learning methodology when it has no access 305 to sellers' private information during training. This can potentially create learning instabilities given 306 that actor-critic training such as A2C generally require that the environment accessed by their critic 307 networks is Markovian (Grondman et al., 2012). Despite this, we find success in this test of "in 308 the wild" learning. The results are displayed in Figure 5.2 and show, despite relaxing this Markov 309 assumption, that the A2C algorithm is able to learn optimal policies for both policy classes (No 310 State and State-based). For No State RL this comes along with lower variance. For State-based 311 RL, the empirical performance is roughly unaffected. We conjecture that the reason behind this 312



Figure 2: Offline learning (left) vs. online, "learning in the wild" (right) performance. The results are averaged over 10 runs, and the shaded regions show 95% conf. intervals.

good performance is related to the class of threshold platform policies. Given a threshold policy, it is possible to predict the overall episode reward based only on the action taken by the policy (the threshold) and ignoring the part of the state that is internal to the sellers (i.e., the Q-matrix and exploration rate). Thus, we see empirically that the additional information that relates to sellers' learning can actually make the platform's learning problem harder.

In the Appendix, we also demonstrate successful experimental results when we replace the use of consumer surplus (1) for reward with a reward that corresponds to the number of agents displayed and the sum of the negated prices offered by sellers. This shows robustness to a possible knowledge gap in knowing the specific functional form of consumer surplus.

322 5.3 Robustness of Learned Platform Rules

As observed in our previous experiments, the Stackelberg-based RL algorithm is effective in learning 323 interventions that maximize consumer surplus for a given economic setting. However, as they are 324 tailored to the economic setting at hand, these interventions can perform poorly when facing settings 325 that are different from those during training. To learn more robust platform rules, we also train 326 with a modified version of the Stackelberg POMDP: at each reward step, with some *random-price* 327 *probability*, sellers quote prices sampled uniformly at random from the price grid. In this way, 328 the platform is rewarded during training for performance that remains robust to prices that are not 329 produced by the Q-learning equilibrium dynamics (given seller costs at training). 330

We evaluate the effect of adding this perturbation-based robustness to the training procedure in settings with different seller costs: in addition to the default c = 1.0, we also test with cost c = 1.38(between the second and the third price in the grid of prices between 0.95 and 2.1) and cost c = 1.67(between the third and the fourth price in the price grid). Here, and for additional realism, we also continue to train the platform rule according to the "in the wild" approach described above, in Section [5.2]

As we see in Figure 5.3 (right), this training approach (and in particular with probability 0.4 of 337 random-price perturbation) succeeds in making the state-based policy much more robust in the 338 face of sellers who experience a different cost environment at test time. The robust, state-based 339 policy displays sellers with higher prices (due to their higher costs), while continuing to significantly 340 mitigate collusion when seller costs are as they were during training. This is also confirmed by the 341 policy visualizations in Figure 5.3 (left), which show how the buy box learned for State-based RL 342 tends to be much more open under this modified training regime. In contrast, the policy learned by 343 No State RL performs very poorly (zero consumer surplus) when tested at costs that differ from those 344 assumed during training, and even under this modified training regime. There is no single threshold 345 that provides a good compromise between performance at cost 1 and handling price perturbations. 346



Figure 3: Left: **Policy visualization**, with number of displayed agents given price selection, averaged over 10 seeds (white–avg. num. sellers displayed 2, black–avg. num. sellers displayed 0). A: No State RL with no price perturbation during training. B: No State RL with 40% random price perturbation(rpp) during training. C: State-based RL with no price perturbation during training. D: State-based RL with 40% random price perturbation during training. Right: **Robustness test**, with buy box policy trained without price perturbation and with price perturbation with prob. 0.4, averaged over 10 runs.

347 6 Conclusion

This work has demonstrated that rules that are effective in preventing collusion by sellers can be 348 learned through a framework that correctly solves the two-level, Stackelberg problem (making use of 349 the platform's commitment power). Specifically, we have introduced the class of *threshold policies* 350 that contain policies that optimize consumer surplus and a learning methodology that is effective in 351 learning optimal leader policies in this class. The interventions we learned are shown to substantially 352 outperform the hand-designed interventions introduced in prior work when the cost environment at 353 test time is as anticipated during training. We also showed how our learned platform interventions 354 can be made more robust when settings are dynamic, with varying seller cost structures, by adopting 355 a suitably-modified training methodology. This also highlights the importance of the state-based 356 platform rule relative to a no-state rule. 357

Interesting future directions include testing our approach in more complex settings, e.g., when sellers' 358 359 costs vary between training episodes. In this case, optimal policy actions are based on the prices quoted during the sellers' equilibration phase, as these prices may provide useful information about 360 the current underlying costs (intuitively, the quoted prices will be higher under higher costs). In 361 this scenario, it may be necessary to represent our platform policies via recurrent neural network, 362 keeping a memory of past prices. Finally, we believe that this approach can also be effective in other 363 applications, e.g., to design and understand effective interventions for the electricity markets studied 364 by Abada and Lambin (2020), a setting where the successful use of RL as a defensive response by a 365 platform is not yet established. 366

367 **References**

Ibrahim Abada and Xavier Lambin. 2020. Artificial Intelligence: Can Seemingly Collusive Outcomes
 Be Avoided? *Available at SSRN 3559308* (2020).

Simon P Anderson and Andre De Palma. 1992. The logit as a model of product differentiation.
 Oxford Economic Papers 44, 1 (1992), 51–67.

³⁷² Enrique Areyan Viqueira, Cyrus Cousins, Yasser Mohammad, and Amy Greenwald. 2019. Empir-

ical Mechanism Design: Designing Mechanisms from Data. In Proceedings of the Thirty-Fifth
 Conference on Uncertainty in Artificial Intelligence (Proceedings of Machine Learning Research,

375 Vol. 115). 1094–1104.

Stephanie Assad, Robert Clark, Daniel Ershov, and Lei Xu. 2020. Algorithmic Pricing and Competi tion: Empirical Evidence from the German Retail Gasoline Market. *CESifo Working Paper Series* 8521 (2020).

Craig Boutilier. 1996. Planning, Learning and Coordination in Multiagent Decision Processes. In
 Proceedings of the 6th Conference on Theoretical Aspects of Rationality and Knowledge. 195–210.

Gianluca Brero, Darshan Chakrabarti, Alon Eden, Matthias Gerstgrasser, Vincent Li, and David
 Parkes. 2021a. Learning Stackelberg Equilibria in Sequential Price Mechanisms. In *ICML Workshop for Reinforcement Learning Theory*.

 Gianluca Brero, Alon Eden, Matthias Gerstgrasser, David C. Parkes, and Duncan Rheingans-Yoo.
 2021b. Reinforcement Learning of Sequential Price Mechanisms. In *Thirty-Fifth AAAI Conference* on Artificial Intelligence, AAAI. 5219–5227.

Erik Brinkman and Michael P. Wellman. 2017. Empirical Mechanism Design for Optimizing Clearing
 Interval in Frequent Call Markets. In *Proceedings of the 2017 ACM Conference on Economics and Computation*. 205–221.

 Benedikt Bünz, Benjamin Lubin, and Sven Seuken. 2018. Designing Core-selecting Payment Rules:
 A Computational Search Approach. In *Proceedings of the 2018 ACM Conference on Economics* and Computation. 109.

Emilio Calvano, Giacomo Calzolari, Vincenzo Denicolò, Joseph E Harrington, and Sergio Pastorello.
 2020b. Protecting consumers from collusive prices due to AI. *Science* 370, 6520 (2020), 1040–
 1042.

Emilio Calvano, Giacomo Calzolari, Vincenzo Denicolo, and Sergio Pastorello. 2020a. Artificial
 intelligence, algorithmic pricing, and collusion. *American Economic Review* 110, 10 (2020),
 3267–97.

C. Cheng, Z. Zhu, B. Xin, and C Chen. 2017. A multi-agent reinforcement learning algorithm based on Stackelberg game. In *6th Data Driven Control and Learning Systems (DDCLS)*. 727–732.

Michael J. Curry, Ping-Yeh Chiang, Tom Goldstein, and John Dickerson. 2020. Certifying Strate gyproof Auction Networks. In *Proc. 33rd Annual Conference on Neural Information Processing Systems*.

Michael J. Curry, Uro Lyi, Tom Goldstein, and John Dickerson. 2022. Learning Revenue-Maximizing
 Auctions With Differentiable Matching. In *Proc. 25th International Conference on Artificial Intelligence and Statistics*. Forthcoming.

Paul Duetting, Zhe Feng, Harikrishna Narasimhan, David C. Parkes, and Sai Srivatsa Ravindranath.
 2019. Optimal Auctions through Deep Learning. In *Proceedings of the 36th International Confer ence on Machine Learning*. 1706–1715.

Tanner Fiez, Benjamin Chasnov, and Lillian J. Ratliff. 2020. Implicit Learning Dynamics in Stackel berg Games: Equilibria Characterization, Convergence Analysis, and Empirical Study. In *Proceed- ings of the 37th International Conference on Machine Learning*. 3133–3144.

⁴¹³ Drew Fudenberg and Jean Tirole. 1991. *Game theory*. MIT press.

Scott Fujimoto, David Meger, Doina Precup, Ofir Nachum, and Shixiang Shane Gu. 2022. Why
 Should I Trust You, Bellman? The Bellman Error is a Poor Replacement for Value Error. *arXiv preprint arXiv:2201.12417* (2022).

⁴¹⁷ Ivo Grondman, Lucian Busoniu, Gabriel AD Lopes, and Robert Babuska. 2012. A survey of actor-⁴¹⁸ critic reinforcement learning: Standard and natural policy gradients. *IEEE Transactions on Systems*,

419 *Man, and Cybernetics, Part C (Applications and Reviews)* 42, 6 (2012), 1291–1307.

Chi Jin, Praneeth Netrapalli, and Michael I. Jordan. 2020. What is Local Optimality in Nonconvex Nonconcave Minimax Optimization?. In *Proceedings of the 37th International Conference on*

⁴²¹ Nonconcave Minimax Optimization?. In *Pi* ⁴²² *Machine Learning*, Vol. 119, 4880–4889.

- ⁴²³ Justin Johnson, Andrew Rhodes, and Matthijs R Wildenbeest. 2020. Platform design when sellers ⁴²⁴ use pricing algorithms. *CEPR Discussion Paper No. DP15504* (2020).
- Patrick R. Jordan, L. Julian Schvartzman, and Michael P. Wellman. 2010. Strategy exploration in
 empirical games. In *9th International Conference on Autonomous Agents and Multiagent Systems*.
 1131–1138.
- Christopher Kiekintveld and Michael P. Wellman. 2008. Selecting strategies using empirical game
 models: an experimental analysis of meta-strategies. In *7th International Joint Conference on Autonomous Agents and Multiagent Systems*. 1095–1101.
- 431 Kevin Kuo, Anthony Ostuni, Elizabeth Horishny, Michael J. Curry, Samuel Dooley, Ping-Yeh Chiang,
- Tom Goldstein, and John P. Dickerson. 2020. ProportionNet: Balancing Fairness and Revenue for
 Auction Design with Deep Learning. *CoRR* abs/2010.06398 (2020).
- Shihui Li, Yi Wu, Xinyue Cui, Honghua Dong, Fei Fang, and Stuart Russell. 2019. Robust Multi Agent Reinforcement Learning via Minimax Deep Deterministic Policy Gradient. *Proceedings of the AAAI Conference on Artificial Intelligence* 33 (July 2019), 4213–4220.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017.
 Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in Neural Information Processing Systems (NeurIPS)* 30 (2017).
- David Mguni, Joel Jennings, Emilio Sison, Sergio Valcarcel Macua, Sofia Ceppi, and Enrique Munoz
 de Cote. 2019. Coordinating the Crowd: Inducing Desirable Equilibria in Non-Cooperative Systems. In Proceedings of the 19th Interpretational Conference on Automatica and MultiAcoust
- tems. In Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent
 Systems. 386–394.
- Rajesh K. Mishra, Deepanshu Vasal, and Sriram Vishwanath. 2020. Model-free Reinforcement
 Learning for Stochastic Stackelberg Security Games. In *59th IEEE Conference on Decision and Control*. 348–353.
- Neehar Peri, Michael J. Curry, Samuel Dooley, and John P. Dickerson. 2021. PreferenceNet: Encoding
 Human Preferences in Auction Design with Deep Learning. In *Proceedings of the 35th Conference on Neural Information Processing Systems*.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
 Dormann. 2021. Stable-Baselines3: Reliable Reinforcement Learning Implementations. *Journal of Machine Learning Research* 22, 268 (2021), 1–8.
- Jad Rahme, Samy Jelassi, Joan Bruna, and S. Matthew Weinberg. 2021b. A Permutation-Equivariant
 Neural Network Architecture For Auction Design. In *Proc. Thirty-Fifth AAAI Conference on Artificial Intelligence*. 5664–5672.
- Jad Rahme, Samy Jelassi, and S. Matthew Weinberg. 2021a. Auction Learning as a Two-Player Game. In *Proc. 9th International Conference on Learning Representations, ICLR*.
- Sailik Sengupta and Subbarao Kambhampati. 2020. Multi-agent Reinforcement Learning in Bayesian
 Stackelberg Markov Games for Adaptive Moving Target Defense. *arXiv:2007.10457 [cs]* (July
 2020). arXiv: 2007.10457.
- Weiran Shen, Binghui Peng, Hanpeng Liu, Michael Zhang, Ruohan Qian, Yan Hong, Zhi Guo,
 Zongyao Ding, Pengjun Lu, and Pingzhong Tang. 2020. Reinforcement mechanism design:
 With applications to dynamic pricing in sponsored search auctions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 2236–2243.
- W. Shen, P. Tang, and S. Zuo. 2019. Automated Mechanism Design via Neural Networks. In
 Proceedings of the 18th International Conference on Autonomous Agents and Multiagent Systems.
- Zhenyu Shi, Runsheng Yu, Xinrun Wang, Rundong Wang, Youzhi Zhang, Hanjiang Lai, and Bo An.
 2020. Learning Expensive Coordination: An Event-Based Deep RL Approach. In *8th International Conference on Learning Representations*.

- Tianmin Shu and Yuandong Tian. 2019. M³RL: Mind-aware Multi-agent Management Reinforce ment Learning. In *7th International Conference on Learning Representations*.
- Andrea Tacchetti, DJ Strouse, Marta Garnelo, Thore Graepel, and Yoram Bachrach. 2019. A Neural
 Architecture for Designing Truthful and Efficient Auctions. *CoRR* 1907.05181 (2019).
- Pingzhong Tang. 2017. Reinforcement mechanism design.. In *Proc. 26th Int Joint Conf. on Art. Intell.* (*IJCAI*). 5146–5150.
- K. Tharakunnel and S. Bhattacharyya. 2007. Leader-follower semi-markov decision problems:
 theoretical framework and approximate solution. In *IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning*. 111–118.
- The Organisation for Economic Co-operation and Development. 2017. Algorithms and Collusion– Note from the European Union. www.oecd.org/competition/algorithms-and-collusion.
- 481 htm

U.S. Federal Trade Commission. 2018. The Competition and Consumer Protection Issues of Algorithms, Artificial Intelligence, and Predictive Analytics, Hearing on Competition and Consumer
 Protection in the 21st Century, U.S. Federal Trade Commission. https://www.ftc.gov/

485 policy/hearings-competition-consumer-protection

- Yevgeniy Vorobeychik, Christopher Kiekintveld, and Michael P. Wellman. 2006. Empirical mech anism design: methods, with application to a supply-chain scenario. In *Proceedings 7th ACM Conference on Electronic Commerce (EC-2006)*. 306–315.
- Kai Wang, Lily Xu, Andrew Perrault, Michael K. Reiter, and Milind Tambe. 2022. Coordinating
 Followers to Reach Better Equilibria: End-to-End Gradient Descent for Stackelberg Games. In
 AAAI Conference on Artificial Intelligence.
- Jakob Weissteiner and Sven Seuken. 2020. Deep Learning—Powered Iterative Combinatorial Auctions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 2284–2293.
- Michael P. Wellman. 2006. Methods for Empirical Game-Theoretic Analysis. In *Proceedings, The Twenty-First National Conference on Artificial Intelligence*. 1552–1556.
- Lily Xu, Andrew Perrault, Fei Fang, Haipeng Chen, and Milind Tambe. 2021. Robust Reinforcement
 Learning Under Minimax Regret for Green Security. In *Proc. 37th Conference on Uncertainty in* Artifical Intelligence (UAI-21).
- Stephan Zheng, Alexander Trott, Sunil Srinivasa, David C Parkes, and Richard Socher. 2022. The
 AI Economist: Optimal Economic Policy Design via Two-level Deep Reinforcement Learning.
 Science Advances 8 (2022).

502 Checklist

- The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or [N/A]. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:
- Did you include the license to the code and datasets? [Yes] See Section ??.
- Did you include the license to the code and datasets? [No] The code and the data are proprietary.
- Did you include the license to the code and datasets? [N/A]

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions

- ⁵¹³ block and only keep the Checklist section heading above along with the questions/answers below.
- 514 1. For all authors...

515 516	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
517	(b) Did you describe the limitations of your work? [Yes] See conclusion.
518 519	(c) Did you discuss any potential negative societal impacts of your work? [N/A] The paper itself is focused on tackling a potential negative societal impact.
520 521	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
522	2. If you are including theoretical results
523	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
524 525	(b) Did you include complete proofs of all theoretical results? [Yes] Proofs will be provided in the appendix.
526	3. If you ran experiments
527 528 529	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes] We will include it in the supplemental material.
530 531	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
532 533	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [Yes] But not in Figure 3 to reduce clutter.
534 535	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the beginning of Section 5.
536	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
537 538	(a) If your work uses existing assets, did you cite the creators? [Yes] We used <i>Stable Baselines3</i> , which we cite at the beginning of Section 5.
539 540	(b) Did you mention the license of the assets? [Yes] We mentioned the <i>Stable Baselines3</i> at the beginning of Section 5.
541	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
542	(d) Did you discuss whether and how consent was obtained from people whose data you're
543	using/curating? [N/A]
544 545	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
546	5. If you used crowdsourcing or conducted research with human subjects
547 548	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
549 550	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
551 552	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]