

Provable Optimization of Quantal Response Leader-Follower Games with Exponentially Large Action Spaces

Jinzhao Li
Purdue University
West Lafayette, IN, USA
li4255@purdue.edu

Carla P. Gomes
Cornell University
Ithaca, NY, USA
gomes@cs.cornell.edu

Daniel Fink, Christopher Wood
Cornell Lab of Ornithology
Ithaca, NY, USA
{df36,clw37}cornell.edu

Yexiang Xue
Purdue University
West Lafayette, IN, USA
yexiang@purdue.edu

ABSTRACT

Leader-follower games involve a leader committing strategies before her followers. We consider quantal response leader-follower games, where the followers' response is probabilistic due to their bounded rationality. Moreover, both the leader's and followers' action spaces are exponentially large with respect to the problem size, hence rendering the overall complexity to solve these games beyond NP-complete. We propose the XOR-Game algorithm, which converges in linear speed towards the equilibrium of convex quantal response leader-follower games (#P-hard to find the equilibrium even though convex). XOR-Game combines stochastic gradient descent with XOR-sampling, a provable sampling approach which transforms highly intractable probabilistic inference into queries to NP oracles. We tested XOR-Game on zero-sum and distribution matching leader-follower games. Experiments show XOR-Game converges faster to a good leader's strategy compared to several baselines. In particular, XOR-Game helps to find the optimal reward allocations for the Avicaching game in the citizen science domain, which harnesses rewards to motivate bird watchers towards tasks of high scientific value.

KEYWORDS

Leader-follower games; Quantal response; XOR-Game

ACM Reference Format:

Jinzhao Li, Daniel Fink, Christopher Wood, Carla P. Gomes, and Yexiang Xue. 2023. Provable Optimization of Quantal Response Leader-Follower Games with Exponentially Large Action Spaces. In *Proc. of the 22nd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2023)*, London, United Kingdom, May 29 – June 2, 2023, IFAAMAS, 19 pages.

1 INTRODUCTION

Leader-follower games, also known as the Stackelberg games [13], involve leaders committing strategies before her followers. Over the years, leader-follower games have been studied extensively with their wide applications in security [47, 56], crowdsourcing [57, 58], AI for social good [19, 41], etc. Recent studies have focused on the participants' *bounded rationality* [46]. In other words, players do

not always play the best moves due to imperfect information or limited computational capacity.

In quantal response leader-follower games, the followers take probabilistic actions due to their bounded rationality. In a Logistic quantal response game, a follower maximizes her utility, albeit together with extreme value distributed latent factors. In the eyes of an observer who does not know the latent factors, the follower's behavior is in an exponential family distribution. Quantal response games have been studied in security games [11], and in games for social good [38]. The authors of [32] proposes an iterative approach to compute a near-optimal strategy for the leader in response to quantal responding adversaries.

In this paper, we consider solving quantal response leader-follower games, in which the sizes of the action spaces of both the leader and the followers grow exponentially quickly w.r.t. the problem size. Exponentially large action spaces are prevalent in real-world games, e.g., in real-time strategy (RTS) games [6] or security games [11]. They pose significant challenges in finding the equilibrium of the game because they prevents the leader from enumerating the entire action space, letting along reasoning about the followers' responses. Although several algorithms have been proposed in computing the equilibrium for normal form games [34, 43] and extensive form games [8, 51, 54], their mathematical programs involve summing over all the followers' actions. This becomes intractable as the number of actions grows exponentially in the size of the game.

We propose XOR-Game, the first algorithm which **converges in linear speed towards the equilibrium of a convex quantal response leader-follower game with exponentially large action spaces**. Despite the game is convex with respect to the leader's strategy, the problem is still at least #P-hard due to the inference of the followers' actions from exponentially many probabilistic choices. Overall XOR-Game optimizes for the leader's objective following a Stochastic Gradient Descent (SGD) process. Our innovation is to harness XOR-sampling in the estimation of the gradient direction of each SGD step towards the optimal leader's strategy. XOR-sampling transforms the highly intractable (#P complete) probabilistic inference and sampling problems into queries to NP oracles using randomly generated XOR constraints. In another view, our XOR-Game algorithm transforms the highly intractable problem of reasoning about the followers' probabilistic actions into problems within the NP complexity class while obtaining provable guarantees on the linear convergence speed and the distances towards

the optimum. Notice other sampling approaches, e.g., MCMC sampling, provide unbiased samples only after an exponential number of burn-in steps. This is impossible in practice, and hence using these sampling approaches cannot result in similar convergence bounds as ours. Our guarantee is also significantly stronger than those offered by e.g., variational approaches [4, 24, 27–29, 45, 55], which are typically lower bounded only and can be arbitrarily loose. Even though the idea of incorporating XOR-sampling appears straightforward, all the theoretic derivation towards linear convergence guarantees cannot borrow from existing theoretic results of SGD. The key difficulty is due to that XOR-sampling only provides constant approximation guarantees for the probability of drawing samples but cannot guarantee unbiased sample estimation, which unfortunately was needed by most prior analysis.

Among many real-world applications, we consider two special cases of convex quantal response leader-follower games. The first is a zero-sum game, where the leader is to minimize the expected utility of the followers. The second is a distribution-matching leader-follower game, which the leader harnesses rewards to move the distribution of the followers’ actions towards a given distribution. Our games have applications in the citizen science domain, where the organizer harnesses rewards to motivate citizen scientists towards tasks with high scientific value. In particular, we apply our game in the recently deployed *Avicaching* game [57, 58] in the *eBird* citizen science framework, where the organizer encourages bird enthusiasts towards bird watching activities in remote and under-sampled locations. The experiment evaluations on both synthetic games as well as on real-world *Avicaching* games show that our XOR-Game is able to produce better leader’s strategies in fewer SGD iterations compared to competing approaches.

2 PRELIMINARIES

2.1 Quantal Response Model

Classic decision theory takes the rational agent assumption, in which agents make perfect choices to maximize their utilities. However, this assumption falls short in the explanation of probabilistic decision-making and occasional deviation from optimal choices. Random utility models were developed to capture the bounded rationality of human decisions [5, 53]. Since its inception, random utility models have been used extensively in modeling human decision-making, ranging from demand prediction [3, 5, 53], behavior modeling [23, 26] to crowd-sourcing [17, 48].

Quantal response games, originated from Quantal choice model [33], were developed from random utility behavior models, and have achieved promising results modeling the bounded rationality of human beings [12, 39]. When faced with $N = 2^n$ choices, where the i -th choice has an observable utility value V_i , quantal response model assumes that the agent’s choice a is to maximize the sum of the utility V_i and a latent factor ϵ_i :

$$a = \arg \max_{i \in \{1, \dots, N=2^n\}} V_i + \epsilon_i. \quad (1)$$

ϵ_i is i.i.d. distributed in the standard Gumbel extreme value distribution Gumbel(0, 1). In other words, the probabilistic nature of the agent’s decision-making is due to the joint optimization of $V_i + \epsilon_i$ rather than V_i only. Noted that ϵ_i is only available to the agent but hidden from the observer. Gumbel noise is well accepted in

literature to account for the stochasticity and/or irrationality of human decision-making [35, 48, 49]. Other types of noises, such as Gaussian noise, lead to other interesting models, such as the probit model. We leave as future work to consider those models. Finally, Gumbel(0,1) is used to simplify the theoretic derivation of our algorithm. Gumbel distributions with other parameters can be considered in a similar way.

Under the random utility model, it can be proven that in the eyes of an observer who do not have access to ϵ_i , the probability that the agent chooses the i -th option is given by the following exponential family distribution:

$$P(i) = \frac{\exp(V_i)}{Z} = \frac{\exp(V_i)}{\sum_{j=1}^{N=2^n} \exp(V_j)}. \quad (2)$$

Here, $Z = \sum_{j=1}^{N=2^n} \exp(V_j)$ is known as the partition function. Various quantities have been calculated for random utility behavior models, including the following expected utility:

THEOREM 2.1. [30] *Under the random utility model, an agent’s expected utility when following the decisions made based on Equation 1 is $\log(\sum_{i=1}^{N=2^n} e^{V_i}) + \gamma$. γ is the Euler-Mascheroni constant.*

Exponential Action Spaces. In this paper, we consider games with exponential many choices. In other words, $N = O(2^n)$ and n denotes the input problem size. For example, in the *Avicaching* game in the citizen science domain where rewards are used to motivate crowdsourcing agents to explore sites with high scientific values, each agent’s choice is represented as a set of locations that the agent explores. Suppose there are n locations, the set of choices are all the sets of locations, the size of which is 2^n .

2.2 Quantal Response Leader-Follower Games

Leader-follower games, also known as Stackelberg games, have attracted much research attention. In a game, the leader commits a strategy before her followers, often resulting in different equilibrium solutions from the Nash equilibrium where both sides commit strategies at the same time.

The random utility behavior model discussed in the previous section leads to quantal response leader-follower games. In a quantal response leader-follower game, the follower’s decision-making follows a random utility model. The equilibrium of a quantal response leader-follower game can be computed through:

$$\begin{aligned} \min_r L(\{V_i, P(i)\}, r), \\ \text{s. t. } a = \arg \max_{i \in \{1, \dots, N=2^n\}} V_i(r) + \epsilon_i, \\ \epsilon_i \stackrel{\text{i.i.d.}}{\sim} \text{Gumbel}(0, 1), \forall i \in \{1, \dots, N = 2^n\}. \end{aligned} \quad (3)$$

Here, the leader’s objective is to minimize L . The follower’s utility function for choice i has part V_i , observable by the leader and depends on the leader’s strategy r . It also contains an extreme value distributed latent part ϵ_i , which is only available to the follower but hidden to the leader.

Compared with the standard leader-follower game, the the quantal response game brings in the probabilistic responses of the followers into consideration, which fits better to reality in many occasions. Due to the latent factor ϵ_i , the follower’s action is probabilistic in the eye of the leader; namely, the follower takes the i -th action

with probability $P(i)$, which has the exponential family form in Equation 2. The objective is written as $L(\{V_i, P(i)\}, r)$ showing that the leader’s objective function can be dependent on his strategy r , the follower’s observable utility V_i and/or the probabilities of making each decision $P(i)$. The formulation in Equation 3 can be used to model both cases where the leader takes pure or mixed strategies. In the pure strategy case, the leader’s action r can be an indicator variable of which action to take. In the mixed strategy case, r becomes a vector listing the probability of taking each action. Since we assume the follower acts according to a quantal response model, she always plays probabilistic (hence mixed) strategies in the eyes of the leader.

2.2.1 Zero-sum Games. While the XOR-Game algorithm can provably optimize many quantal response leader-follower games, we mainly consider two variants for this paper. The first variant we consider is the zero sum case, where the leader is to minimize the expected utility of the follower. In other words, the leader’s objective is (according to Theorem 2.1):

$$L_0(\{V_i, P(i)\}, r) = \log \left(\sum_{i=1}^{N=2^n} e^{V_i(r)} \right) + \gamma. \quad (4)$$

We consider a special case where V_i is linear in r , namely, $V_i(r) = \theta_i^T r + \phi_i$. Here θ_i measures the influence of rewards on the leader’s action. ϕ_i represents the intrinsic utility in choosing action i and can vary across actions (even though it does not depend on r). We show that the game in this case is convex in r (proof is in the supplementary materials):

THEOREM 2.2. *When $V_i(r) = \theta_i^T r + \phi_i$, the zero-sum quantal response leader-follower game is convex in r . Moreover, the gradient $\nabla L_0(r)$ has the following form of an expectation:*

$$\nabla L_0(r) = \mathbb{E}_{P(i)} [\theta_i] = \sum_{i=1}^{N=2^n} P(i) \theta_i. \quad (5)$$

It is known in a zero-sum matrix game, the Stackelberg equilibrium matches exactly to the Nash equilibrium [52, 61]. Nevertheless, we would like to point out that in our definition of a quantal response game, the follower’s action is designed to maximize her utility function V_i in addition to an unobserved ϵ_i (Equation 1), where V_i depends on the complete information of the leader’s strategy r . This definition implicitly assumes the leader commits her strategy before the follower. The Nash equilibrium, in this setting, can be difficult to be properly defined.

2.2.2 Distribution Matching Games. Another variant we consider in this paper is the game where the leader would like to stimulate certain behaviors from the follower. In particular, the leader would like to match the probability distribution of the follower’s actions P to a desired distribution Q . The difference between two distributions is measured by Kullback–Leibler (KL) divergence. In other words, the leader’s objective is:

$$L_{DM}(\{V_i, P(i)\}, r) = KL(Q||P) = \sum_{i=1}^{N=2^n} Q(i) \log \left(\frac{Q(i)}{P(i)} \right). \quad (6)$$

This formulation is specially useful for mechanism design problems, e.g. in [20, 40, 44]. In cases where certain follower’s actions increase

the social welfare, the leader would set high Q values to promote these actions. Again, we focus on the special case where V_i is linear in r in this paper:

THEOREM 2.3. *When $V_i(r) = \theta_i^T r + \phi_i$ is linear in r , the distribution matching leader-follower game is convex in r and the gradient $\nabla L_{DM}(r)$ can be represented as:*

$$\nabla L_{DM} = \mathbb{E}_{P(i)} [\theta_i] - \mathbb{E}_{Q(i)} [\theta_i] = \sum_{i=1}^{N=2^n} P(i) \theta_i - \sum_{i=1}^{N=2^n} Q(i) \theta_i. \quad (7)$$

This theorem’s proof is in the supplementary materials.

Avicaching Game in Citizen Science. As a specific example, we consider a leader-follower game in the citizen science domain, where the leader (citizen science game organizer) harnesses limited rewards to encourage citizen scientists (followers) to conduct observations in remote and undersampled locations. We look into the Avicaching game hosted in the eBird citizen science program [57], where rewards are used to encourage bird watchers (citizen scientists, or followers) to visit undersampled Avicaching sites, which have large scientific values but are inconvenient and/or less interesting to visit than traditional hotspots. Each bird watcher’s choice is characterized by a set of locations $L \subseteq \{l_1, \dots, l_n\}$, which represents the set of spots one plan to visit during a bird watching trip. The organizer, in this case, harnesses reward $r = (r_1, \dots, r_n)$ to stimulate visits to under-sampled locations. Here, r_i is the reward that a bird watcher receives when he visits location l_i . The bird watchers’ (followers’) utility $V_L(r)$ models both the intrinsic utilities to visit the location set L as well as the reward received.

2.3 XOR Sampling

Sampling from a combinatorial space has a formal complexity of #P-complete, the difficulty of which is beyond NP-completeness. Luckily, the recently proposed XOR-Sampling algorithm, as the result of a rich line of works using streamlining randomized constraints [1, 2, 9, 15, 16, 21, 22], provides a constant approximation guarantee on the probabilities of the samples generated. XOR-sampling transforms the highly intractable sampling problem into queries to NP-oracles while obtaining provable guarantees.

The high-level idea of XOR-sampling is to harness randomized constraints to guarantee the randomness of the samples generated. Consider a simple case where $w(x) : \{0, 1\}^n \rightarrow \{0, 1\}$ is a binary function and one would like to obtain a sample from the solution space $\{x : w(x) = 1\}$ uniformly at random. Querying a NP oracle returns one x satisfying $w(x) = 1$ albeit not at random. XOR-sampling works by querying NP oracles to find x which satisfies $w(x) = 1$ and subject to a few randomized XOR constraints. It can be proven that each additional XOR constraint removes approximately half of the solutions to $w(x) = 1$ at random. Hence, once a desirable number of XOR constraints are added and the resulting space has only one solution, it can be proven that the only solution remaining is a random one from the original space $\{x : w(x) = 1\}$. In this way, XOR-sampling is able to bound the probability of obtaining each sample within a constant multiplicative factor of its ground-truth probability. For general weighted functions XOR-sampling has similar guarantees, although the sampling process becomes

more complex. Our proposed XOR-Game algorithm depends on the following approximation bounds:

THEOREM 2.4. (Ermon et al., 2013)¹ Let $1 < \delta < \sqrt{2}$, $0 < \gamma < 1$, $w : \{0, 1\}^n \rightarrow \mathbb{R}^+$ be an unnormalized weight function. $P(x) \propto w(x)$ is the normalized distribution. Then, with probability at least $1 - \gamma$, XOR-Sampling(w, δ, γ) succeeds and outputs a sample x_0 using $O(-n \log(1 - 1/\sqrt{\delta}) \log(-n/\gamma \log(1 - 1/\sqrt{\delta})))$ NP-oracle queries. Upon success, each x_0 is produced with probability $P'(x_0)$. We have

$$\frac{1}{\delta} P(x_0) \leq P'(x_0) \leq \delta P(x_0).$$

Moreover, let $\phi : \{0, 1\}^n \rightarrow \mathbb{R}$ be a function mapping binary vectors to \mathbb{R} . Denote $\phi(x)^+ = \max\{\phi(x), 0\}$ and $\phi(x)^- = \min\{\phi(x), 0\}$ as the positive and negative part of $\phi(x)$. Then the expectation of one sampled $\phi(x)$ satisfies,

$$\frac{1}{\delta} \mathbb{E}_{P(x)}[\phi(x)^+] \leq \mathbb{E}_{P'(x)}[\phi(x)^+] \leq \delta \mathbb{E}_{P(x)}[\phi(x)^+],$$

$$\delta \mathbb{E}_{P(x)}[\phi(x)^-] \leq \mathbb{E}_{P'(x)}[\phi(x)^-] \leq \frac{1}{\delta} \mathbb{E}_{P(x)}[\phi(x)^-].$$

Algorithm 1: XOR-Game₀

Input : $r_0, \{\theta_i\}_{i=1}^N, \{\phi_i\}_{i=1}^N$
Params : $T, K, \eta, \delta, \gamma$

for $t = 0$ **to** T **do**

$k \leftarrow 1$;

while $k \leq K$ **do**

$P(i) \propto \exp(\theta_i^T r_t + \phi_i)$;

$l' \leftarrow \text{XOR-Sampling}(P(i), \delta, \gamma)$;

if $l' \neq \text{Failure}$ **then**

$l'_k \leftarrow l'$;

$k \leftarrow k + 1$;

end

end

$\bar{g}_t \leftarrow \frac{1}{K} \sum_{k=1}^K \theta_{l'_k}$;

$r_{t+1} = r_t - \eta \bar{g}_t$;

end

return $\bar{r}_T = \frac{1}{T} \sum_{t=1}^T r_t$

3 XOR-GAME

The challenge in solving quantal response leader follower games is the intractable probabilistic inference over the follower's strategies. In this paper, we consider games in which the follower's action space is exponentially large. In other words, the number of actions N is of size $O(2^n)$, where n is the problem size. These games are prevalent in real world. See the Avicaching game presented in the experiment section for an example. Notice we assume there are a compact representation for all θ_i 's and ϕ_i 's. Even though there are $2N$ vectors of these, we assume the availability of efficient functions $\theta(i)$ and $\phi(i)$. When given i , they return θ_i and ϕ_i , respectively. The length of encoding both functions $\theta(i)$ and $\phi(i)$ are within $O(n)$, i.e., the length of the problem description. In this setup, the quantal

¹The details of the discretization scheme and the choices of parameters of the original algorithm which yield the bound of this form is discussed in [14].

Algorithm 2: XOR-Game_{DM}

Input : $r_0, \{Q(i)\}_{i=1}^N, \{\theta_i\}_{i=1}^N, \{\phi_i\}_{i=1}^N$
Params : $T, K, S, \eta, \delta, \gamma$

$j \leftarrow 1$;

while $j \leq S$ **do**

$l'' \leftarrow \text{XOR-Sampling}(Q, \delta, \gamma)$;

if $l'' \neq \text{Failure}$ **then**

$l''_j \leftarrow l''$; $j \leftarrow j + 1$;

end

end

for $t = 0$ **to** T **do**

$k \leftarrow 1$;

while $k \leq K$ **do**

$P(i) \propto \exp(\theta_i^T r_t + \phi_i)$;

$l' \leftarrow \text{XOR-Sampling}(P(i), \delta, \gamma)$;

if $l' \neq \text{Failure}$ **then**

$l'_k \leftarrow l'$; $k \leftarrow k + 1$;

end

end

$\bar{g}_t \leftarrow \frac{1}{K} \sum_{k=1}^K \theta_{l'_k} - \frac{1}{S} \sum_{j=1}^S \theta_{l''_j}$;

$r_{t+1} = r_t - \eta \bar{g}_t$;

end

return $\bar{r}_T = \frac{1}{T} \sum_{t=1}^T r_t$

response leader-follower game is at least #P-hard, even limiting to the convex games considered in Theorem 2.2 and 2.3. This is because it is already #P-hard to compute the partition function in $P(i)$ in Equation 2. In other words, it is already #P-complete to evaluate the leader's objective function even for a fixed strategy.

We propose XOR-Game, which **converges towards the equilibrium of convex quantal response leader-follower games in linear number of stochastic gradient descent iterations**. The XOR-Game algorithm should enjoy the convergence bound for a wide variety of quantal response games. However, the actual algorithms and the convergence bounds slightly differ across different game setups, due to differences in estimating the derivatives and their correspondingly different approximation bounds given by XOR-sampling. In this paper, we demonstrate such convergence bounds on the aforementioned zero-sum and distribution matching quantal response leader-follower games. However, we are confident that similar guarantees generalize to many other games.

The algorithm variants for solving zero sum game and distribution matching game are shown in Algorithm 1 and Algorithm 2. The procedures of XOR-Game₀ and XOR-Game_{DM} have minimal differences. Both algorithms apply SGD to find the optimal reward r that minimizes the leader's objective. r_0 is the initialization of the reward vector. The follower's observable utility is $V_i(r) = \theta_i^T r + \phi_i$. Samples generated from XOR-sampling are used to estimate the expectations in the gradient calculation (according to Equation 5 and 7). Because XOR-Sampling has a failure rate, repeated sampling is used until a desired number of samples are obtained. K samples are drawn from the behavior model of followers, and S

samples are from the targeting model. XOR-Sampling takes parameters (δ, γ) . After the gradient estimation, r_{t+1} from the next iteration moves from r_t following the negative gradient direction. η is the step size of SGD. Finally after T SGD steps, the average of r_1, \dots, r_T is returned as the output. Denote the total variance $\text{Var}_{P(i)}(\theta_i) = \mathbb{E}_{P(i)}(\|\theta_i\|_2^2) - \|\mathbb{E}_{P(i)}(\theta_i)\|_2^2$. We can show that the convergence bound for XOR-Game₀ to solve the zero sum game is:

THEOREM 3.1. (Convergence for zero-sum game) *In a zero sum quantal response leader follower game with objective in Equation 4 and $V_i(r) = \theta_i^T r + \phi_i$, r^* attains the minimum of the leader's objective. \bar{r}_T is the output of XOR-Game₀ starting from r_0 and running T SGD iterations. In iteration t of SGD, g_t is the estimated gradient, i.e., $r_{t+1} = r_t - \eta g_t$. If $\max_P \text{Var}_{P(i)}(\theta_i) \leq \sigma^2$, $\|r_t - r^*\|_2 \leq R$, $\eta \leq (2 - \delta^2)/(\sigma^2 \delta)$, $\max_P \|\mathbb{E}_{P(i)}(\theta_i^+)\|_2 \leq G$, and $\max_P \|\mathbb{E}_{P(i)}(\theta_i^+)\|_2 \leq G$, where $\theta^+ = \max\{\theta, 0\}$ and $\theta^- = \min\{\theta, 0\}$, we have*

$$\mathbb{E}[L_0(\bar{r}_T)] - L_0(r^*) \leq \frac{\delta \|r_0 - r^*\|_2^2}{2\eta T} + \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta\delta(\delta^2 - 1)G^2 + 2(\delta^2 - 1)GR.$$

Theorem 3.1 proves that XOR-Game₀ converges to the equilibrium of the zero-sum game in a linear number of SGD steps in addition to a few vanishing terms. Here, the first term on the right-hand side scales inversely proportional to the number of SGD iterations T , suggesting a linear convergence speed towards the equilibrium. The second term can be reduced by increasing K , the number of XOR samples in estimating $\mathbb{E}_{P(i)}[\theta_i]$. The third and the fourth terms are the products of constants with $(\delta^2 - 1)$, which can be minimized with a more accurate, yet more time-consuming XOR-sampling (bringing δ closer to 1). The proof of Theorem 3.1 shares the same high-level idea with Theorem 3.4. The detailed proof is left to the supplementary materials. The convergence bound for XOR-Game_{DM}, the algorithm variant to solve the distribution matching game depends on a stronger condition:

Definition 3.2. (Match signs at every dimension) A group of vectors $\Theta = \{\theta_1, \dots, \theta_N\}$ matches signs at every dimension, if for any two vectors $\theta_i, \theta_j \in \Theta$, $\theta_i = (\theta_{i1}, \dots, \theta_{iL})^T$, $\theta_j = (\theta_{j1}, \dots, \theta_{jL})^T$, for any dimension $k \in \{1, \dots, L\}$, we have $\theta_{ik}\theta_{jk} \geq 0$.

The provable guarantee for XOR-Game_{DM} requires all θ_i in the distribution matching game to match signs at every dimension. This requirement is not too stringent. As we have pointed out, distribution matching leader follower games are usually seen in mechanism design problems, where the leader searches for a strategy to maximize certain behaviors from the followers. Here, the leader's strategy r typically represents the incentives offered to the follower. θ_i in this case becomes indicator variables whether certain incentives are earned if the follower takes action i . Due to this reason, all θ_i are non-negative, satisfying the matching signs condition. With these definitions, the convergence bound for the distribution matching leader follower game is as follows:

THEOREM 3.3. (Convergence for distribution matching game) *Suppose a distribution matching leader-follower game has the objective in Equation 6. $V_i(r) = \theta_i^T r + \phi_i$. Denote r^* as the optimal leader's strategy. \bar{r}_T is the output of the XOR-Game_{DM}. Suppose $\{\theta_1, \dots, \theta_N\}$ match signs at every dimension, $\max_P \text{Var}_{P(i)}(\theta_i) \leq$*

σ^2 , $\max_P \|\mathbb{E}_{P(i)}(\theta_i)\|_2 \leq G$, $\text{Var}_{Q(i)}(\theta_i) \leq \sigma^2$, $\|\mathbb{E}_{Q(i)}(\theta_i)\|_2 \leq G$, when $1 < \delta < \sqrt{2}$ is used in XOR-sampling and the SGD step size $\eta \leq (2 - \delta^2)/(\sigma^2 \delta)$, $\|r_t - r^*\|_2 \leq R$ for all r_1, \dots, r_T , we have:

$$\mathbb{E}[L_{DM}(\bar{r}_T)] - L_{DM}(r^*) \leq \frac{\delta \|r_0 - r^*\|_2^2}{2\eta T} + (\delta^2 - 1) \left[\sqrt{2}GR + 2\eta \left(\frac{\sigma^2 + G^2}{\min\{K, S\}} + \delta G^2 \right) \right] + 2\eta(\delta^2 + 1) \frac{\sigma^2 + G^2}{\min\{K, S\}}. \quad (8)$$

To interpret this inequality, the first term on the right-hand side of inequality 8 scales inversely proportional to T , showing a linear convergence speed towards the optimal leader's strategy r^* . The second term is the product of $(\delta^2 - 1)$ and a constant (all terms in the square bracket). This term can be minimized with more accurate (yet more expensive) XOR-sampling, bringing in δ closer to 1. The term in the second line scales inversely proportional to $\min\{K, S\}$, which can be minimized by increasing K and S ; e.g., drawing more samples. In summary, this theorem still shows a linear convergence bound and two tails terms which can be minimized via better sampling. The proof of Theorem 3.3 depends on the following Theorem 3.4, which was motivated by Theorem 3 in [14]. Nevertheless, Theorem 3 in [14] does not apply to the case where the difference of two XOR sampling processes are used to estimate the gradient. We therefore need to come up with novel proof techniques for distribution matching games, which yields the following Theorem 3.4:

THEOREM 3.4. *Suppose function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth convex. $r^* = \arg \min_r f(r)$. At any point r , the gradient $\nabla f(r)$ can be decomposed into $\nabla f(r) = \nabla p(r) - \nabla q(r)$. At the t -th iteration of SGD, $g_t = k_t - l_t$ is the estimated gradient, i.e., $r_{t+1} = r_t - \eta g_t$. k_t (or l_t) is the estimation of $\nabla p(r_t)$ (or $\nabla q(r_t)$). $\{k_t, l_t, \nabla p(r_t), \nabla q(r_t)\}$ match signs at every dimension. If $\text{Var}(k_t) \leq \sigma^2$, $\text{Var}(l_t) \leq \sigma^2$, $\|\mathbb{E}[k_t]\|_2^2 \leq G^2$, $\|\mathbb{E}[l_t]\|_2^2 \leq G^2$, and there exists $1 < c < \sqrt{2}$, s.t.*

$$\begin{aligned} \frac{1}{c} [\nabla p(r_t)]^+ &\leq \mathbb{E}[k_t^+] \leq c [\nabla p(r_t)]^+ \\ c [\nabla p(r_t)]^- &\leq \mathbb{E}[k_t^-] \leq \frac{1}{c} [\nabla p(r_t)]^- \\ \frac{1}{c} [\nabla q(r_t)]^+ &\leq \mathbb{E}[l_t^+] \leq c [\nabla q(r_t)]^+ \\ c [\nabla q(r_t)]^- &\leq \mathbb{E}[l_t^-] \leq \frac{1}{c} [\nabla q(r_t)]^-. \end{aligned}$$

Let $R = \max_t \|r_t - r^*\|$, with $\eta \leq \frac{2-c^2}{Lc}$, $\bar{r}_T = \frac{1}{T} \sum_{t=1}^T r_t$, we have:

$$\mathbb{E}[f(\bar{r}_T)] - f(r^*) \leq \frac{c}{2\eta T} \|r_0 - r^*\|_2^2 + \left(c - \frac{1}{c}\right) \left(\sqrt{2}GR + 2\eta(\sigma^2 + G^2)\right) + 2\eta \left(c + \frac{1}{c}\right) \sigma^2. \quad (9)$$

The proof of Theorem 3.3 is to apply Theorem 3.4 to the objective function of XOR-Game_{DM}. Notice that L_{DM} is L -smooth when the total variation $\max_P \text{Var}_{P(i)}(\theta_i)$ is bounded (proved in a lemma). Those 4 constraints on the expectation of estimated gradients can be achieved by tuning parameters of XOR-Sampling.

To prove Theorem 3.4, we need the following lemmas. The proofs of these lemmas are left in the supplementary materials:

LEMMA 3.5. *Suppose f is convex. $r^* = \arg \min_r f(r)$. At the t -th iteration of SGD, $g_t = k_t - l_t$ is the estimated gradient. $\{k_t, l_t, \nabla p(r_t),$*

$\nabla q(r_t)$ match signs at every dimension, and there exists $1 < c < \sqrt{2}$, s.t. $\frac{1}{c}[\nabla p(r_t)]^+ \leq \mathbb{E}[k_t^+] \leq c[\nabla p(r_t)]^+$, $c[\nabla p(r_t)]^- \leq \mathbb{E}[k_t^-] \leq \frac{1}{c}[\nabla p(r_t)]^-$, $\frac{1}{c}[\nabla q(r_t)]^+ \leq \mathbb{E}[l_t^+] \leq c[\nabla q(r_t)]^+$, $c[\nabla q(r_t)]^- \leq \mathbb{E}[l_t^-] \leq \frac{1}{c}[\nabla q(r_t)]^-$, we have:

$$\langle \nabla p(r_t), \mathbb{E}[k_t] \rangle \geq \frac{1}{c} \|\mathbb{E}[k_t]\|_2^2, \quad (10)$$

$$\langle \nabla q(r_t), \mathbb{E}[l_t] \rangle \geq \frac{1}{c} \|\mathbb{E}[l_t]\|_2^2, \quad (11)$$

$$\langle \nabla p(r_t), \mathbb{E}[l_t] \rangle \leq c \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle, \quad (12)$$

$$\langle \nabla q(r_t), \mathbb{E}[k_t] \rangle \leq c \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle. \quad (13)$$

LEMMA 3.6. Suppose all variables and pre-conditions are defined as in Theorem 3.4. In particular, $\text{Var}(k_t) \leq \sigma^2$, $\text{Var}(l_t) \leq \sigma^2$, $\|\mathbb{E}[k_t]\|_2^2 \leq G^2$, $\|\mathbb{E}[l_t]\|_2^2 \leq G^2$, we have

$$\| \text{Tr}[\text{Cov}(k_t, l_t)] \| = |\mathbb{E}\langle k_t, l_t \rangle - \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle| \leq \sigma^2, \quad (14)$$

$$\mathbb{E}\langle k_t, l_t \rangle \leq \sigma^2 + G^2. \quad (15)$$

LEMMA 3.7. Suppose all variables and conditions are defined as in Theorem 3.4. We have:

$$\begin{aligned} \langle \nabla f(r_t), r_t - r^* \rangle &\leq c \langle \mathbb{E}[k_t] - \mathbb{E}[l_t], r_t - r^* \rangle + \\ &\quad \sqrt{2} \left(c - \frac{1}{c} \right) GR. \end{aligned} \quad (16)$$

PROOF. (Formal proof of Theorem 3.4) By L-smoothness of f , for the t -th iteration,

$$\begin{aligned} f(r_{t+1}) &\leq f(r_t) + \langle \nabla f(r_t), r_{t+1} - r_t \rangle + \frac{L}{2} \|r_{t+1} - r_t\|_2^2 \\ &= f(r_t) - \eta \langle \nabla p(r_t) - \nabla q(r_t), k_t - l_t \rangle + \frac{L\eta^2}{2} \|k_t - l_t\|_2^2 \\ &= f(r_t) + \frac{L\eta^2}{2} \|k_t - l_t\|_2^2 - \\ &\quad \eta \langle \langle \nabla p(r_t), k_t \rangle - \langle \nabla q(r_t), k_t \rangle - \langle \nabla p(r_t), l_t \rangle + \langle \nabla q(r_t), l_t \rangle \rangle. \end{aligned}$$

Take the expectation w.r.t. k_t and l_t on both sides, and notice Equations 10, 11, 12, 13 in Lemma 3.5, we have:

$$\begin{aligned} \mathbb{E}[f(r_{t+1})] &\leq f(r_t) - \frac{\eta}{c} [\|\mathbb{E}[k_t]\|_2^2 + \|\mathbb{E}[l_t]\|_2^2] + \\ &\quad 2\eta c \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle + \frac{L\eta^2}{2} \mathbb{E}[\|k_t - l_t\|_2^2]. \end{aligned}$$

Notice $\text{Var}(k_t) = \mathbb{E}[\|k_t\|_2^2] - \|\mathbb{E}[k_t]\|_2^2$, $\text{Var}(l_t) = \mathbb{E}[\|l_t\|_2^2] - \|\mathbb{E}[l_t]\|_2^2$, and $\text{Tr}[\text{Cov}(k_t, l_t)] = \mathbb{E}\langle k_t, l_t \rangle - \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle$, we further rewrite the right-hand side as:

$$\begin{aligned} &\mathbb{E}[f(r_{t+1})] \\ &\leq f(r_t) - \frac{\eta}{c} [\mathbb{E}[\|k_t\|_2^2] - \text{Var}(k_t) + \mathbb{E}[\|l_t\|_2^2] - \text{Var}(l_t)] + \\ &\quad 2\eta c [\mathbb{E}\langle k_t, l_t \rangle - \text{Tr}[\text{Cov}(k_t, l_t)]] + \frac{L\eta^2}{2} \mathbb{E}[\|k_t - l_t\|_2^2]. \end{aligned}$$

After re-arranging terms, the right-hand side again becomes:

$$\mathbb{E}[f(r_{t+1})] \leq f(r_t) - \frac{\eta(2 - L\eta c)}{2c} \mathbb{E}[\|k_t - l_t\|_2^2] + \text{tail}. \quad (17)$$

$$\begin{aligned} \text{tail} &= \frac{\eta}{c} (\text{Var}(k_t) + \text{Var}(l_t)) - 2\eta c \text{Tr}[\text{Cov}(k_t, l_t)] + \\ &\quad 2\eta \left(c - \frac{1}{c} \right) \mathbb{E}\langle k_t, l_t \rangle. \end{aligned}$$

Using $\text{Var}(k_t) \leq \sigma^2$, $\text{Var}(l_t) \leq \sigma^2$ and Lemma 3.6, we have:

$$\text{tail} \leq 2\eta \left(c - \frac{1}{c} \right) (\sigma^2 + G^2) + 2\eta \left(c + \frac{1}{c} \right) \sigma^2. \quad (18)$$

For simplicity, denote the right-hand side of the previous inequality as a constant C_1 . Hence, Equation 17 becomes:

$$\mathbb{E}[f(r_{t+1})] \leq f(r_t) - \frac{\eta(2 - L\eta c)}{2c} \mathbb{E}[\|k_t - l_t\|_2^2] + C_1.$$

Using $\eta \leq \frac{2-c^2}{Lc}$, we can further simplify this inequality to:

$$\mathbb{E}[f(r_{t+1})] \leq f(r_t) - \frac{\eta c}{2} \mathbb{E}[\|k_t - l_t\|_2^2] + C_1.$$

Because f is convex, $f(r_t) \leq f(r^*) + \langle \nabla f(r_t), r_t - r^* \rangle$. Follow the previous inequality we get:

$$\mathbb{E}[f(r_{t+1})] \leq f(r^*) + \langle \nabla f(r_t), r_t - r^* \rangle - \frac{\eta c}{2} \mathbb{E}[\|k_t - l_t\|_2^2] + C_1.$$

Because of Lemma 3.7, we can further rewrite the previous inequality as

$$\begin{aligned} \mathbb{E}[f(r_{t+1})] &\leq f(r^*) + c \langle \mathbb{E}[k_t] - \mathbb{E}[l_t], r_t - r^* \rangle + \\ &\quad \sqrt{2} \left(c - \frac{1}{c} \right) GR - \frac{\eta c}{2} \mathbb{E}[\|k_t - l_t\|_2^2] + C_1. \end{aligned}$$

Define $C_2 = C_1 + \sqrt{2} \left(c - \frac{1}{c} \right) GR$, and notice $k_t - l_t = g_t$, we can write

$$\begin{aligned} \mathbb{E}[f(r_{t+1})] &\leq f(r^*) + c \langle \mathbb{E}[g_t], r_t - r^* \rangle - \frac{\eta c}{2} \mathbb{E}[\|g_t\|_2^2] + C_2 \\ &= f(r^*) + \frac{c}{2\eta} \left(2\eta \langle \mathbb{E}[g_t], r_t - r^* \rangle - \eta^2 \mathbb{E}[\|g_t\|_2^2] \right) + C_2 \\ &= f(r^*) + \frac{c}{2\eta} \mathbb{E} \left[2\eta \langle g_t, r_t - r^* \rangle - \eta^2 \|g_t\|_2^2 \right] + C_2. \end{aligned}$$

From the second last to the last equation, we also take expectation w.r.t. r_t on both sides. The equality holds because the randomness of g_t come from the sampling step at the t -th iteration, which is independent of r_t (whose randomness come from the first $t-1$ iterations). Because $r_{t+1} = r_t - \eta g_t$, we have $\|r_t - r^*\|_2^2 - \|r_{t+1} - r^*\|_2^2 = 2\eta \langle g_t, r_t - r^* \rangle - \eta^2 \|g_t\|_2^2$. Hence we have:

$$\mathbb{E}[f(r_{t+1})] \leq f(r^*) + \frac{c}{2\eta} \mathbb{E}[\|r_t - r^*\|_2^2 - \|r_{t+1} - r^*\|_2^2] + C_2. \quad (19)$$

By summing up Equation 19 for $t = 0, \dots, T-1$, we get

$$\sum_{t=0}^{T-1} \mathbb{E}[f(r_{t+1}) - f(r^*)] \leq \frac{c\|r_0 - r^*\|_2^2}{2\eta} + TC_2. \quad (20)$$

Finally, by Jensen's inequality, $Tf(\bar{r}_T) \leq \sum_{t=1}^T f(r_t)$,

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbb{E}[f(r_{t+1}) - f(r^*)] &= \mathbb{E} \left[\sum_{t=1}^T f(r_t) \right] - Tf(r^*) \\ &\geq T\mathbb{E}[f(\bar{r}_T)] - Tf(r^*). \end{aligned} \quad (21)$$

Combining equations 20, 21, we have

$$\mathbb{E}[f(\bar{r}_T)] \leq f(r^*) + \frac{c\|r_0 - r^*\|_2^2}{2\eta T} + C_2, \quad (22)$$

which is exactly the equation in Theorem 3.4. \square

To quantify the computational complexity of XOR-Game, we prove the following theorem in the supplementary materials detailing the number of queries to NP oracles needed for XOR-Game_{DM} and XOR-Game₀. The proof of this Theorem is again left in the supplementary materials.

THEOREM 3.8. *XOR-Game₀ in Algorithm 1 uses $O(-Tn \log(1 - 1/\sqrt{\delta}) \log(-n/\gamma \log(1 - 1/\sqrt{\delta})) + TK)$ queries to NP oracles. XOR-Game_{DM} in Algorithm 2 uses $O(-Tn \log(1 - 1/\sqrt{\delta}) \log(n/\gamma \log(1 - 1/\sqrt{\delta})) + TK + S)$ queries to NP oracles.*

4 EXPERIMENTS

We demonstrate empirical evidence that XOR-Game outperforms a few competing approaches in the speed and the quality of the solutions found for both the quantal response zero-sum leader-follower games and distribution-matching games. Our evaluation is conducted on a synthetic benchmark set and a behavior model learned from real-world data obtained from the Avicaching game, which promotes bird watchers to collect data in remote and undersampled locations using the so-called Avicaching points [57, 58]. The baseline approaches we consider are: (1) BRQR algorithm [59], which minimizes the leader’s objective in quantal response stackelberg games. Their approach is based on a full gradient descend (GD) optimizer, hence needs to go over all the follower’s actions in each iteration and is only applicable on games with small action spaces. (2) gibbs_game, which uses SGD to minimize the leader’s objective but utilizes Gibbs sampling in the estimation of the gradient direction. (3) bp_game, which uses samples generated from the marginal probabilities computed via loopy belief propagation during SGD, [31, 37, 60] and (4) cbp_game, which harnesses the recently proposed BP chain method in generating samples in SGD [18]. For the fairness of comparisons, the leader’s objective L_0 for the zero-sum game is computed using an exact model counter Ace [10]. The leader’s objective L_{DM} for the distribution matching game is the KL-divergence, which is computed using Ace and XOR-sampling. The estimated KL-divergence is close to the groundtruth due to the constant approximation guarantee of XOR-sampling and the exactness of Ace. Additional details are in the supplementary materials.

In synthetic and real-world experiments, we use the Avicaching game as the background. In the Avicaching game, the leader (Avicaching game organizer) harnesses rewards to motivate the followers (bird watchers) to visit remote and under-sampled locations. The rewards are in the form of virtual Avicaching points, which marks the participants’ contributions to science. At the end of each season of the Avicaching game, a lottery is drawn from which Avicaching participants have opportunities to win birding gears based on how many Avicaching points they have contributed. In both the synthetic and the real-world experiments, one action that one Avicaching participant can take is to visit a set of locations L . In practice, we assume bird watchers only choose between locations historically documented in the eBird dataset [50] hence we have information for all the locations. We assume the probability that one Avicaching participant visit a set of locations L is given by:

$$P(L) \propto \exp(w_r \theta_L^T r + w_f FL + L^T WL). \quad (23)$$

Here, we use $P(L)$ instead of $P(i)$ because each action is characterized by a set of locations. We assume L is represented as a vector

Table 1: Comparison between XOR-Game₀ and BRQR

Size	Loss _{XOR}	Loss _{BRQR}	Time _{XOR}	Time _{BRQR}
2 ²	0.0425	0.0071	147.59s	4.10s
2 ⁴	0.0677	0.0149	154.96s	28.11s
2 ⁸	0.0346	0.0139	196.48s	46.16s
2 ¹⁶	0.0338	0.0178	302.24s	499.57s
2 ³²	0.0814	N/A	699.36s	>3h
2 ⁶⁴	0.0799	N/A	8476.04s	>3h

of indicator variables of visited locations. $w_f FL + L^T WL$ is represented using symbol ϕ during theoretic derivation. Each column of F includes features associated with each location, such as its landscape composition, proximity to water, etc, which affects bird watchers’ intrinsic utilities in visiting these locations. W is a matrix characterizing the changing of utilities when visiting multiple locations (e.g., bird watchers typically do not prefer visiting multiple locations of the same type). θ_L is a vector of indicator variables of whether visiting location set L receives each reward. w_r and w_f are the relative importance of rewards and location features. For distribution matching game, we assume $Q(L)$ has the same form as $P(L)$ although with different parameters to promote visits to remote and under-sampled locations.

Validation on Small Games. We first validate that XOR-Game finds close-to-optimal leader’s strategies on small sized games. Specifically, we focus on the zero-sum quantal response game. BRQR is used as the baseline for comparison. The difference in leader’s utility values between the equilibrium and the ones found by the algorithms are shown as Loss_{XOR} and Loss_{BRQR} in Table 1. Here “Size” represents the number of different location sets a follower can visit, Time_{XOR} and Time_{BRQR} are the running times of different algorithms respectively. We can see from the table that both XOR-Game and BRQR find close-to-optimal leader’s strategies. Initially XOR-Game takes longer to converge, but BRQR cannot scale to modest sized games as it runs out of a 3-hour time limit for a game with action space of $\geq 2^{32}$. XOR-Game still produces near optimal solutions in this size. Further details of this experiment (in particular, the speed the two algorithms converge to these solutions) are left to the supplementary materials.

Evaluation on Large Synthetic Benchmarks. We further evaluate the performance of XOR-Game on both the zero-sum game and the distribution matching game on large synthetic benchmarks. In these experiments, we intentionally increase the dimensionality of the reward vector r to be quadratic in the number of locations to increase the difficulty of benchmarks. To be specific, we let $\theta_L = \text{vector}(LL^T) = (l_1 l_1, l_1 l_2, \dots, l_1 l_n, \dots, l_n l_n)^T$ where L is the location set vector $L = (l_1, \dots, l_n)^T$. This is to assume one participant can receive a unique reward r_{ij} by visiting location pair (i, j) . We enforce each r_{ij} to be non-negative and no greater than 1. During SGD, when r_{ij} becomes negative (or bigger than 1), we reset it to be 0 (or 1). Additional details are in the supplementary materials.

Figure 1 (left and middle) shows the performance of various algorithms as the optimization progresses. Here, each curve shows the leader’s objective averaged over 20 benchmarks. For each benchmark, we let all algorithms start from the same initial solution. When computing the average, we normalize the objective function

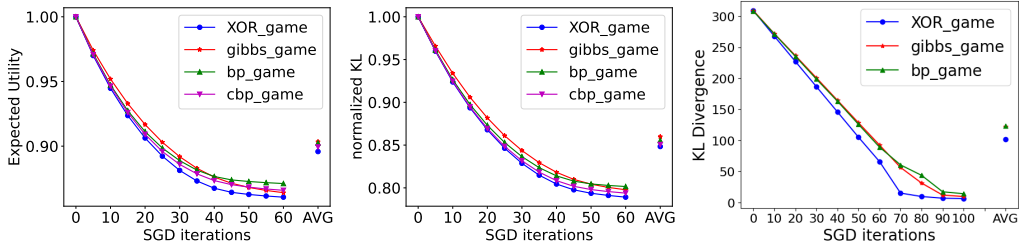


Figure 1: (Left and Middle) XOR-Game converges faster and to better solutions compared with competing approaches on synthetic datasets. (Left) Experiment on the zero sum game. (Middle) Experiment on the distribution matching game. (Right) XOR-Game converges faster than competing approaches on a behavior model learned from data collected from a real-world Avicaching game. X-axis shows the number of SGD iterations. Y-axis shows the leader’s objective function $L(L_0$ or $L_{DM})$. AVG depicts $L(\bar{r}_T)$.

values against that of the initial solution so each algorithm always starts from an objective function value of 1.

Notably, XOR-Game descends to the best solutions within the least number of SGD iterations for both the zero-sum games and the distribution matching games. We would like to point out that XOR-sampling in this case is efficient in obtaining the samples, even though XOR-sampling has to answer NP-complete queries. In particular, it roughly takes 1 second for XOR-sampling to obtain 100 samples during SGD, but in general it takes 50 seconds for Gibbs sampling (300 MCMC steps), 2.8 seconds for belief propagation, and 4700 seconds for chained belief propagation (cbp). Because cbp is so slow, we use 100 samples in the gradient estimation for all other approaches but only 10 samples for cbp.

Evaluation on the Avicaching Game. We then evaluate all approaches on a behavior model learned from real-world data collected from the Avicaching Game. The data comes from an actual field deployment of the Avicaching game in the eBird crowdsourcing platform between March 27 and October 29, 2015 (30 weeks) in Tompkins and Cortland counties of the New York State. A set of 50 Avicaching locations were selected, which were all publicly accessible but received no visits prior to the game. The goal of the Avicaching game was to shift the bird watchers’ efforts from traditional bird watching hot spots to these Avicaching locations, harnessing Avicaching points. The numbers of Avicaching points offered for each visit to these Avicaching locations were updated every Monday. The Avicaching game was remarkably effective during this field deployment. It was reported in [57] that 19% of the bird watching effort in these two counties were shifted to these under-sampled Avicaching locations.

Before evaluating algorithms, we first learn a behavior model in the form of Equation 23 for all eBird participants. Because the field study gives an independent reward to each Avicaching location, we set θ_L to be L . $\theta_L^T r = L^T r$ represents the total reward from visiting the location set L . F includes landscape features obtained from the 2011 National Land Cover Database (NLCD). $L^T W L$ represents the change in utility functions for visiting multiple locations. Overall, $w_f F L + L^T W L$ represents the intrinsic utility of visiting locations L . Each data point consists of the set of locations L one bird watcher visits and the corresponding reward r of the week. Parameters w_r , w_f , and W are learned using Contrastive Divergence [7].

We run various algorithms for the distribution matching game to minimize the KL-divergence between the learned probability density $P(L)$ with a manually designed $Q(L)$, which promotes the visiting of under-sampled Avicaching locations and suppresses the visiting to others. The rewards were set to be greater than 0 but less than 100 for each location (same order of magnitude as the actual field deployment). Additional details in terms of learning $P(L)$ and $Q(L)$ can be found in the supplementary materials.

Figure 1(right) demonstrates that XOR-Game descends to an optimal reward allocation faster than competing approaches. All benchmarks start with identically initialized rewards. We manually inspected the solutions. The final solutions of all approaches reach almost zero for the KL-divergence, suggesting a possibility to match the learned probability density $P(L)$ to $Q(L)$ using available rewards. Nevertheless, we cannot conclude that the Avicaching game participants will act according to Q if we had the opportunities to deploy the rewards into the field. This is because all calculations are based on a learned behavior model from historical data. We cannot guarantee how much the learned model captures the subtle aspects of human decision-making and new visiting patterns may emerge as the human behavior changes with the introduction of the Avicaching game. On average, the wall-clock time for each method is: XOR_game(24h), bp_game(21h), and gibbs_game(10h). The cbp_game is excluded for comparison because it takes >8h per SGD iteration. In summary, XOR-Game requires the least number of SGD iterations to descend to the best leader’s strategies among all benchmark algorithms. XOR-game completes in a reasonable time, has good empirical performance and provable guarantees.

5 CONCLUSION

We proposed XOR-Game to solve the convex quantal response leader-follower games with exponentially large action spaces. XOR-Game has a linear convergence speed towards the equilibrium of the leader-follower games. Our approach is based on an integration of XOR-Sampling and stochastic gradient descent, transforming the otherwise #P-hard problem into queries within the NP complexity class, while obtaining guarantees for the convergence speed. The experiments on both synthetic and real-world Avicaching games show that XOR-Game outperforms other baseline methods and hence prove its great potential for real-world applications.

ACKNOWLEDGEMENTS

This research was supported by NSF grants IIS-1850243, CCF-1918327. We thank the anonymous reviewers for their comments and suggestions.

REFERENCES

- [1] Dimitris Achlioptas, Zayd Hammoudeh, and Panos Theodoropoulos. 2018. Fast and Flexible Probabilistic Model Counting. In *International Conference on Theory and Applications of Satisfiability Testing*.
- [2] Dimitris Achlioptas and Panos Theodoropoulos. 2017. Probabilistic model counting with short XORs. In *International Conference on Theory and Applications of Satisfiability Testing*. Springer, 3–19.
- [3] Victor Aguirregabiria and Pedro Mira. 2010. Dynamic discrete choice structural models: A survey. *Journal of Econometrics* 156, 1 (2010), 38–67. <https://doi.org/10.1016/j.jeconom.2009.09.007>
- [4] David Belanger and Andrew McCallum. 2016. Structured prediction energy networks. In *International Conference on Machine Learning*. 983–992.
- [5] Moshe Ben-Akiva and Steven R. Lerman. 1985. *Discrete choice analysis: theory and application to travel demand*. MIT Press.
- [6] Michael Buro. 2003. Real-time strategy games: A new AI research challenge. In *IJCAI*. Vol. 2003. 1534–1535.
- [7] Miguel Á. Carreira-Perpiñán and Geoffrey E. Hinton. 2005. On Contrastive Divergence Learning. In *AISTATS*.
- [8] Jakub Cerný, Viliam Lisý, Branislav Bošanský, and Bo An. 2021. Computing Quantal Stackelberg Equilibrium in Extensive-Form Games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 5260–5268.
- [9] Supratik Chakraborty, Daniel J. Fremont, Kuldeep S. Meel, Sanjit A. Seshia, and Moshe Y. Vardi. 2014. Distribution-aware Sampling and Weighted Model Counting for SAT. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence (AAAI)*.
- [10] Mark Chavira, Adnan Darwiche, and Manfred Jaeger. 2006. Compiling relational Bayesian networks for exact inference. *Int. J. Approx. Reasoning* (2006).
- [11] Kam-Fung Cheung and Michael GH Bell. 2021. Attacker-defender model against quantal response adversaries for cyber security in logistics management: An introductory study. *European Journal of Operational Research* 291, 2 (2021), 471–481.
- [12] Renaud Chicoisne and Fernando Ordóñez. 2016. Risk averse Stackelberg security games with quantal response. In *International Conference on Decision and Game Theory for Security*. Springer, 83–100.
- [13] Vincent Conitzer and Tuomas Sandholm. 2006. Computing the optimal strategy to commit to. In *EC '06*.
- [14] Fan Ding, Jianzhu Ma, Jinbo Xu, and Yexiang Xue. 2021. XOR-CD: Linearly Convergent Constrained Structure Generation. In *Proceedings of the Thirty-eighth International Conference on Machine Learning (ICML)*.
- [15] Fan Ding and Yexiang Xue. 2021. XOR-SGD: Provable Convex Stochastic Optimization for Decision-making under Uncertainty. In *Proceedings of the 37th Conference on Uncertainty in Artificial Intelligence (UAI)*.
- [16] Stefano Ermon, Carla P. Gomes, Ashish Sabharwal, and Bart Selman. 2013. Embed and Project: Discrete Sampling with Universal Hashing. In *Advances in Neural Information Processing Systems (NIPS)*.
- [17] Stefano Ermon, Yexiang Xue, Russell Toth, Bistra N. Dilkina, Richard Bernstein, Theodoros Damoulas, Patrick Clark, Steve DeGloria, Andrew Mude, Christopher Barrett, and Carla P. Gomes. 2015. Learning Large-Scale Dynamic Discrete Choice Models of Spatio-Temporal Preferences with Application to Migratory Pastoralism in East Africa. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI)*.
- [18] Ding Fan and Yexiang Xue. 2020. Contrastive Divergence Learning with Chained Belief Propagation. In *International Conference on Probabilistic Graphical Models*.
- [19] Fei Fang, Peter Stone, and Milind Tambe. 2015. When Security Games Go Green: Designing Defender Strategies to Prevent Poaching and Illegal Fishing. In *IJCAI*.
- [20] Noah Golowich, Hari Krishna Narasimhan, and David C Parkes. 2018. Deep Learning for Multi-Facility Location Mechanism Design.. In *IJCAI*. 261–267.
- [21] Carla P. Gomes, Ashish Sabharwal, and Bart Selman. 2006. Model Counting: A New Strategy for Obtaining Good Bounds. In *Proceedings of the 21st National Conference on Artificial Intelligence*.
- [22] Carla P. Gomes, Ashish Sabharwal, and Bart Selman. 2007. Near-Uniform Sampling of Combinatorial Spaces Using XOR Constraints. In *Advances in Neural Information Processing Systems*.
- [23] Vassilis Argyrou Hajivassiliou. 1991. Simulation Estimation Methods for Limited Dependent Variable Models. *Handbook of Statistics* 11 (1991), 519–543.
- [24] T. Hazan and A. Shashua. 2010. Norm-Product Belief Propagation: Primal-Dual Message-Passing for Approximate Inference. *Information Theory, IEEE Transactions on* 56, 12 (Dec 2010), 6294–6316. <https://doi.org/10.1109/TIT.2010.2079014>
- [25] Tom Heskes and Kees Albers. 2003. Approximate inference and constrained optimization. In *In 19th UAI*. 313–320.
- [26] Joel L. Horowitz. 1993. SEMIPARAMETRIC ESTIMATION OF A WORK-TRIP MODE-CHOICE MODEL. *Journal of Econometrics* 58 (1993), 49–70.
- [27] Bai Jiang, Tung-Yu Wu, Yifan Jin, Wing H Wong, et al. 2018. Convergence of contrastive divergence algorithm in exponential family. *The Annals of Statistics* 46, 6A (2018), 3067–3098.
- [28] Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional Variational Autoencoder with Adversarial Learning for End-to-End Text-to-Speech. *arXiv preprint arXiv:2106.06103* (2021).
- [29] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [30] Xiao Bo Li. 2018. *Convex Optimization and Online Learning: Their Applications in Discrete Choice Modeling and Pricing*. ProQuest Dissertations Publishing.
- [31] Qiang Liu and Alexander Ihler. 2012. Belief Propagation for Structured Decision Making. In *Uncertainty in Artificial Intelligence (UAI)*. AUAI Press, Corvallis, Oregon, 523–532.
- [32] Xiao Bo Ma, Bo An, Mengchen Zhao, Xiapu Luo, Lei Xue, Zhenhua Li, Tony TN Miu, and Xiaohong Guan. 2019. Randomized security patrolling for link flooding attack detection. *IEEE Transactions on Dependable and Secure Computing* 17, 4 (2019), 795–812.
- [33] Daniel L. McFadden. 1976. Quantal choice analysis: A survey. *Annals of Economic and Social Measurement, Volume 5, number 4* (1976), 363–390.
- [34] Richard D McKelvey, Andrew M McLennan, and Theodore L Turocy. 2006. Gambit: Software tools for game theory. (2006).
- [35] Sanjog Misra and Sanjib K Mohanty. 2008. Estimating bargaining games in distribution channels. In *Working Paper*. Citeseer.
- [36] Joris M. Mooij. 2010. libDAI: A Free and Open Source C++ Library for Discrete Approximate Inference in Graphical Models. *Journal of Machine Learning Research* 11 (Aug. 2010), 2169–2173.
- [37] Kevin P Murphy, Yair Weiss, and Michael I Jordan. 1999. Loopy belief propagation for approximate inference: An empirical study. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 467–475.
- [38] Theo Offerman, Arthur Schram, and Joep Sonnemans. 1998. Quantal response models in step-level public good games. *European Journal of Political Economy* 14, 1 (1998), 89–100.
- [39] Özlem Ömer. 2018. Dynamics of the us housing market: A quantal response statistical equilibrium approach. *Entropy* 20, 11 (2018), 831.
- [40] Balázs Pejó and Gergely Biczók. 2020. Corona Games: Masks, Social Distancing and Mechanism Design. In *Proceedings of the 1st ACM SIGSPATIAL International Workshop on Modeling and Understanding the Spread of COVID-19 (Seattle, WA, USA) (COVID-19)*. Association for Computing Machinery, New York, NY, USA, 24–31. <https://doi.org/10.1145/3423459.3430757>
- [41] A. Perrault, Bryan Wilder, Eric Ewing, Aditya Mate, Bistra N. Dilkina, and Milind Tambe. 2020. End-to-End Game-Focused Learning of Adversary Behavior in Security Games. In *AAAI*.
- [42] Wei Ping and Alex Ihler. 2017. Belief propagation in conditional RBMs for structured prediction. In *Artificial Intelligence and Statistics*. PMLR, 1141–1149.
- [43] Sainath Sanga, Venkata Srimad Siddhardh Nadendla, and Sajal K Das. 2021. Maximizing Social Welfare in Selfish Multi-Modal Routing using Strategic Information Design for Quantal Response Travelers. *arXiv preprint arXiv:2111.15069* (2021).
- [44] Avi Segal, Kobi Gal, Ece Kamar, and Eric Horvitz. 2018. Optimizing Interventions via Offline Policy Evaluation: Studies in Citizen Science. In *AAAI-18: Thirty-Second AAAI Conference on Artificial Intelligence*.
- [45] Huajie Shao, Shuochao Yao, Dachun Sun, Aston Zhang, Shengzhong Liu, Dongxin Liu, Jun Wang, and Tarek Abdelzaher. 2020. Controlvae: Controllable variational autoencoder. In *International Conference on Machine Learning*. PMLR, 8655–8664.
- [46] Herbert A. Simon. 1990. *Bounded Rationality*. Palgrave Macmillan UK, London, 15–18. https://doi.org/10.1007/978-1-349-20568-4_5
- [47] Arunesh Sinha, Fei Fang, Bo An, Christopher Kiekintveld, and Milind Tambe. 2018. Stackelberg Security Games: Looking Beyond a Decade of Success. In *IJCAI*.
- [48] Hossein Azari Soufiani, David C. Parkes, and Lirong Xia. 2012. Random Utility Theory for Social Choice. In *NIPS*.
- [49] Karunakaran Sudhir. 2001. Structural analysis of manufacturer pricing in the presence of a strategic retailer. *Marketing Science* 20, 3 (2001), 244–264.
- [50] Brian L. Sullivan, Jocelyn L. Aycrigg, Jessie H. Barry, Rick E. Bonney, Nicholas Bruns, Caren B. Cooper, Theo Damoulas, André A. Dhondt, Tom Dieterich, Andrew Farnsworth, Daniel Fink, John W. Fitzpatrick, Thomas Fredericks, Jeff Gerbracht, Carla Gomes, Wesley M. Hochachka, Marshall J. Iliff, Carl Lagoze, Frank A. La Sorte, Matthew Merrifield, Will Morris, Tina B. Phillips, Mark Reynolds, Amanda D. Rodewald, Kenneth V. Rosenberg, Nancy M. Trautmann, Andrea Wiggins, David W. Winkler, Weng-Keen Wong, Christopher L. Wood, Jun Yu, and Steve Kelling. 2014. The eBird enterprise: An integrated approach to development and application of citizen science. *Biological Conservation* 169, 0 (2014), 31–40.
- [51] Czander Tan. 2022. Double-Oracle Deep Reinforcement Learning for Handling Exponential Action Space in Sequential Stackelberg Security Games. (2022).
- [52] Yasuhiro Tanaka. 2018. Stackelberg type dynamic zero-sum game with leader and follower.

- [53] K. Train, D. McFadden, and M. Ben-Akiva. 1987. The Demand for Local Telephone Service: A Fully Discrete Model of Residential Call Patterns and Service Choice. *RAND Journal of Economics* 18 (1987).
- [54] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575, 7782 (2019), 350–354.
- [55] Martin J. Wainwright, Tommi S. Jaakkola, and Alan S. Willsky. 2003. Tree-reweighted belief propagation algorithms and approximate ML estimation by pseudo-moment matching. In *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*.
- [56] Haifeng Xu, Zinovi Rabinovich, Shaddin Dughmi, and Milind Tambe. 2015. Exploring Information Asymmetry in Two-Stage Security Games. In *AAAI*.
- [57] Yexiang Xue, Ian Davies, Daniel Fink, Christopher Wood, and Carla P. Gomes. 2016. Avicaching: A Two Stage Game for Bias Reduction in Citizen Science. In *Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*.
- [58] Yexiang Xue, Ian Davies, Daniel Fink, Christopher Wood, and Carla P. Gomes. 2016. Behavior Identification in Two-stage Games for Incentivizing Citizen Science Exploration. In *Proceedings of the 22nd International Conference on Principles and Practice of Constraint Programming (CP)*.
- [59] Rong Yang, Christopher Kiekintveld, Fernando Ordóñez, Milind Tambe, and Richard John. 2013. Improving resource allocation strategies against human adversaries in security games: An extended study. *Artificial Intelligence* 195 (2013), 440–469.
- [60] Jonathan S. Yedidia, William T. Freeman, and Yair Weiss. 2003. Exploring Artificial Intelligence in the New Millennium. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, Chapter Understanding Belief Propagation and Its Generalizations, 239–269. <http://dl.acm.org/citation.cfm?id=779343.779352>
- [61] Zhengyu Yin, Dmytro Korzhyk, Christopher Kiekintveld, Vincent Conitzer, and Milind Tambe. 2010. Stackelberg vs. Nash in security games: interchangeability, equivalence, and uniqueness. In *AAMAS 2010*.

A ALGORITHM AND PROOFS

A.1 Proofs for Theorems in Section 2

PROOF. (Theorem 2.2) First, consider the derivative of the leader's objective function of $L_0(r)$ in Equation 4 w.r.t. r ,

$$\begin{aligned}\nabla L_0(r) &= \nabla_r \log \left(\sum_{i=1}^N e^{\theta_i^T r + \phi_i} \right) \\ &= \frac{1}{\sum_{j=1}^N e^{\theta_j^T r + \phi_j}} \sum_{i=1}^N e^{\theta_i^T r + \phi_i} \theta_i \\ &= \sum_{i=1}^n P(i) \theta_i = \mathbb{E}_{P(i)} [\theta_i]\end{aligned}$$

Then compute the Hessian matrix $H(L_0)$,

$$\begin{aligned}H(L_0) &= \nabla(\nabla L_0(r)) \\ &= \sum_{i=1}^N P(i) \theta_i \theta_i^T - \left(\sum_{i=1}^N P(i) \theta_i \right) \left(\sum_{j=1}^N P(j) \theta_j \right)^T \\ &= \text{Cov}(\theta_i, \theta_j).\end{aligned}$$

Apparently, the Hessian matrix of the leader's objective function $L_0(r)$ is in the form of a co-variance matrix which is positive semi-definite. Thus the convexity of $L_0(r)$ is proved. \square

PROOF. (Theorem 2.3) Taking derivative of the objective function $L_{DM}(r)$ w.r.t. r ,

$$\begin{aligned}\nabla L_{DM} &= \nabla_r \left(\sum_{i=1}^N Q(i) \log \frac{Q(i)}{P(i)} \right) \\ &= - \nabla_r \left(\sum_{i=1}^N Q(i) \log P(i) \right) \\ &= \sum_{i=1}^N \frac{Q(i)}{P(i)} \frac{\left(\sum_{j=1}^N e^{\theta_j^T r + \phi_j} \theta_j \right) e^{\theta_i^T r + \phi_i}}{\left(\sum_{j=1}^N e^{\theta_j^T r + \phi_j} \right)^2} - \sum_{i=1}^N \frac{Q(i)}{P(i)} \frac{e^{\theta_i^T r + \phi_i} \theta_i}{\sum_{j=1}^N e^{\theta_j^T r + \phi_j}} \\ &= \sum_{i=1}^N Q(i) \sum_{j=1}^N P(j) \theta_j - \sum_{i=1}^N Q(i) \theta_i \\ &= \mathbb{E}_{P(i)} [\theta_i] - \mathbb{E}_{Q(i)} [\theta_i]\end{aligned}$$

The Hessian matrix $H(L_{DM})$ is

$$\begin{aligned}H(L_{DM}) &= \nabla(\nabla L_{DM}) = \nabla(\mathbb{E}_{P(i)} [\theta_i]) \\ &= \sum_{i=1}^N P(i) \theta_i \theta_i^T - \left(\sum_{i=1}^N P(i) \theta_i \right) \left(\sum_{j=1}^N P(j) \theta_j \right)^T \\ &= \text{Cov}(\theta_i, \theta_j).\end{aligned}$$

The Hessian matrix $H(L_{DM})$ is in the form of a co-variance matrix which is positive semi-definite. So L_{DM} is convex w.r.t. r . \square

A.2 Proofs for Zero-sum Quantal Response Leader Follower Games

The XOR-Game algorithm for solving zero-sum quantal response leader follower game is shown in Algorithm 1. The performance of this algorithm is guaranteed by Theorem 3.1. Our definition for

$[\nabla L_0(r)]^+$ (or $[\nabla L_0(r)]^-$) in the proofs to the zero-sum game are as follows:

$$\begin{aligned}[\nabla L_0(r)]^+ &= \mathbb{E}_{P(i)} [\theta_i^+] = \sum_{i=1}^N P(i) [\theta_i^+], \\ [\nabla L_0(r)]^- &= \mathbb{E}_{P(i)} [\theta_i^-] = \sum_{i=1}^N P(i) [\theta_i^-].\end{aligned}$$

Here, $[f]^+ = \max\{f, 0\}$ extracts the positive part of f , and $[f]^- = \min\{f, 0\}$ extracts the negative part of f . Notice in the previous definition, we first extract the positive (negative) part of each θ_i then take the expectation, the result of which can be different from first taking the expectation then extracting the positive (negative) part. It is straightforward to see that

$$\nabla L_0(r) = [\nabla L_0(r)]^+ + [\nabla L_0(r)]^-.$$

In the lemmas below we sometimes use f to represent L_0 , in which case $[\nabla f(r_t)]^+$ ($[\nabla f(r_t)]^-$) assume that we first extract the positive (negative) part then take the expectation as well. This definition carries over to the distribution matching leader follower games, although there we impose the condition that all $\{\theta_1, \dots, \theta_N\}$ match signs at every dimension. Under that condition, the order of taking the expectation and extracting the positive (negative) parts do not affect the final result.

LEMMA A.1. Suppose $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex. $r^* = \arg \min_r f(r)$. At the t -th iteration of SGD, g_t is the estimated gradient in Algorithm 1: $r_{t+1} = r_t - \eta g_t$. Suppose $\|\mathbb{E}[g_t^+]\|_2 \leq G$, $\|\mathbb{E}[g_t^-]\|_2 \leq G$, $\|r_t - r^*\|_2 \leq R$, we have:

$$|\langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^-] \rangle| \leq G^2, \quad (24)$$

$$|\langle \mathbb{E}[g_t^+], [r - r^*]^- \rangle| \leq GR, \quad (25)$$

$$|\langle \mathbb{E}[g_t^-], [r - r^*]^+ \rangle| \leq GR. \quad (26)$$

PROOF. (Lemma A.1) Use Cauchy-Schwarz Inequality, we have,

$$\begin{aligned}|\langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^-] \rangle| &\leq \|\mathbb{E}[g_t^+]\|_2 \|\mathbb{E}[g_t^-]\|_2 \leq G^2, \\ |\langle \mathbb{E}[g_t^+], [r - r^*]^- \rangle| &\leq \|\mathbb{E}[g_t^+]\|_2 \|[r - r^*]^- \|_2 \\ &= \|\mathbb{E}[g_t^+]\|_2 \|\min\{r - r^*, \mathbf{0}\}\|_2 \\ &\leq \|\mathbb{E}[g_t^+]\|_2 \|r - r^*\|_2 \\ &\leq GR \\ |\langle \mathbb{E}[g_t^-], [r - r^*]^+ \rangle| &\leq \|\mathbb{E}[g_t^-]\|_2 \|[r - r^*]^+ \|_2 \\ &= \|\mathbb{E}[g_t^-]\|_2 \|\max\{r - r^*, \mathbf{0}\}\|_2 \\ &\leq \|\mathbb{E}[g_t^-]\|_2 \|r - r^*\|_2 \\ &\leq GR\end{aligned}$$

This completes the proof. \square

LEMMA A.2. Suppose $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex. $r^* = \arg \min_r f(r)$. At the t -th iteration of SGD, g_t is the estimated gradient in Algorithm 1: $r_{t+1} = r_t - \eta g_t$. Suppose $\|\mathbb{E}[g_t^+]\|_2 \leq G$, $\|\mathbb{E}[g_t^-]\|_2 \leq G$, $\|r_t - r^*\|_2 \leq R$. If there exists $1 < c < \sqrt{2}$, s.t. $\frac{1}{c} [\nabla f(r_t)]^+ \leq \mathbb{E}[g_t^+] \leq c [\nabla f(r_t)]^+$, $c [\nabla f(r_t)]^- \leq \mathbb{E}[g_t^-] \leq \frac{1}{c} [\nabla f(r_t)]^-$, then we

have:

$$\frac{1}{c} \|\mathbb{E}[g_t]\|_2^2 \leq \langle \nabla f(r_t), \mathbb{E}[g_t] \rangle + 2(c - \frac{1}{c})G^2, \quad (27)$$

$$\langle \nabla f(r_t), r_t - r^* \rangle \leq c \langle \mathbb{E}[g_t], r_t - r^* \rangle + 2(c - \frac{1}{c})GR. \quad (28)$$

PROOF. (Lemma A.2) Because

$$\begin{aligned} \frac{1}{c} [\nabla f(r_t)]^+ &\leq \mathbb{E}[g_t^+] \leq c [\nabla f(r_t)]^+, \\ c [\nabla f(r_t)]^- &\leq \mathbb{E}[g_t^-] \leq \frac{1}{c} [\nabla f(r_t)]^-, \end{aligned}$$

We have:

$$\begin{aligned} \frac{1}{c} \|\mathbb{E}[g_t^+]\|_2^2 &= \frac{1}{c} \langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^+] \rangle \leq \langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^+] \rangle \\ &\leq c \langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^+] \rangle = c \|\mathbb{E}[g_t^+]\|_2^2. \end{aligned} \quad (29)$$

Similarly,

$$\begin{aligned} \frac{1}{c} \|\mathbb{E}[g_t^-]\|_2^2 &= \frac{1}{c} \langle \mathbb{E}[g_t^-], \mathbb{E}[g_t^-] \rangle \leq \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^-] \rangle \\ &\leq c \langle \mathbb{E}[g_t^-], \mathbb{E}[g_t^-] \rangle = c \|\mathbb{E}[g_t^-]\|_2^2. \end{aligned} \quad (30)$$

For cross terms, we have:

$$\langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^-] \rangle \geq c \langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^-] \rangle, \quad (31)$$

$$\langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^+] \rangle \geq c \langle \mathbb{E}[g_t^-], \mathbb{E}[g_t^+] \rangle. \quad (32)$$

Notice that

$$\begin{aligned} \frac{1}{c} \|\mathbb{E}[g_t]\|_2^2 &= \frac{1}{c} \|\mathbb{E}[g_t^+] + \mathbb{E}[g_t^-]\|_2^2 \\ &= \frac{1}{c} \left(\|\mathbb{E}[g_t^+]\|_2^2 + \|\mathbb{E}[g_t^-]\|_2^2 + 2 \langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^-] \rangle \right) \end{aligned}$$

Use the results in Equation 29, 30, 31, 32, we can further derive:

$$\begin{aligned} \frac{1}{c} \|\mathbb{E}[g_t]\|_2^2 &\leq \langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^+] \rangle + \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^-] \rangle + \\ &\quad \frac{1}{c^2} \langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^-] \rangle + \frac{1}{c^2} \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^+] \rangle \\ &= \langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^+] \rangle + \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^-] \rangle + \\ &\quad \langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^-] \rangle + \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^+] \rangle + \\ &\quad \left(\frac{1}{c^2} - 1 \right) \left(\langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^-] \rangle + \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^+] \rangle \right) \\ &= \langle \nabla f(r_t), \mathbb{E}[g_t] \rangle + \left(\frac{1}{c^2} - 1 \right) \langle [\nabla f(r_t)]^+, \mathbb{E}[g_t^-] \rangle + \\ &\quad \left(\frac{1}{c^2} - 1 \right) \langle [\nabla f(r_t)]^-, \mathbb{E}[g_t^+] \rangle \\ &\leq \langle \nabla f(r_t), \mathbb{E}[g_t] \rangle + \left(\frac{1}{c} - c \right) \langle \mathbb{E}[g_t^+], \mathbb{E}[g_t^-] \rangle + \\ &\quad \left(\frac{1}{c} - c \right) \langle \mathbb{E}[g_t^-], \mathbb{E}[g_t^+] \rangle. \end{aligned}$$

Applying Equation 24 from Lemma A.1 to the last equation, we get Equation 27.

To prove Equation 28, first notice:

$$\begin{aligned} \frac{1}{c} \langle \mathbb{E}[g_t^+], [r_t - r^*]^+ \rangle &\leq \langle [\nabla f(r_t)]^+, [r_t - r^*]^+ \rangle \\ &\leq c \langle \mathbb{E}[g_t^+], [r_t - r^*]^+ \rangle, \\ \frac{1}{c} \langle \mathbb{E}[g_t^-], [r_t - r^*]^- \rangle &\leq \langle [\nabla f(r_t)]^-, [r_t - r^*]^- \rangle \\ &\leq c \langle \mathbb{E}[g_t^-], [r_t - r^*]^- \rangle, \\ c \langle \mathbb{E}[g_t^+], [r_t - r^*]^- \rangle &\leq \langle [\nabla f(r_t)]^+, [r_t - r^*]^- \rangle \\ &\leq \frac{1}{c} \langle \mathbb{E}[g_t^+], [r_t - r^*]^- \rangle, \\ c \langle \mathbb{E}[g_t^-], [r_t - r^*]^+ \rangle &\leq \langle [\nabla f(r_t)]^-, [r_t - r^*]^+ \rangle \\ &\leq \frac{1}{c} \langle \mathbb{E}[g_t^-], [r_t - r^*]^+ \rangle. \end{aligned}$$

Then we have:

$$\begin{aligned} &\langle \nabla f(r_t), r_t - r^* \rangle \\ &= \langle [\nabla f(r_t)]^+ + [\nabla f(r_t)]^-, [r_t - r^*]^+ + [r_t - r^*]^- \rangle \\ &= \langle [\nabla f(r_t)]^+, [r_t - r^*]^+ \rangle + \langle [\nabla f(r_t)]^+, [r_t - r^*]^- \rangle + \\ &\quad \langle [\nabla f(r_t)]^-, [r_t - r^*]^+ \rangle + \langle [\nabla f(r_t)]^-, [r_t - r^*]^- \rangle \\ &\leq c \langle \mathbb{E}[g_t^+], [r_t - r^*]^+ \rangle + c \langle \mathbb{E}[g_t^-], [r_t - r^*]^- \rangle + \\ &\quad \frac{1}{c} \langle \mathbb{E}[g_t^-], [r_t - r^*]^+ \rangle + \frac{1}{c} \langle \mathbb{E}[g_t^+], [r_t - r^*]^- \rangle \\ &= c \langle \mathbb{E}[g_t^+], [r_t - r^*]^+ \rangle + c \langle \mathbb{E}[g_t^-], [r_t - r^*]^- \rangle + \\ &\quad c \langle \mathbb{E}[g_t^-], [r_t - r^*]^+ \rangle + c \langle \mathbb{E}[g_t^+], [r_t - r^*]^- \rangle + \\ &\quad \left(\frac{1}{c} - c \right) \left(\langle \mathbb{E}[g_t^-], [r_t - r^*]^+ \rangle + \langle \mathbb{E}[g_t^+], [r_t - r^*]^- \rangle \right) \\ &= \left(\frac{1}{c} - c \right) \left(\langle \mathbb{E}[g_t^-], [r_t - r^*]^+ \rangle + \langle \mathbb{E}[g_t^+], [r_t - r^*]^- \rangle \right) + \\ &\quad c \langle \mathbb{E}[g_t], [r_t - r^*] \rangle \end{aligned}$$

Using Equation 25, 26, we have the previous line of equation cannot be greater than $c \langle \mathbb{E}[g_t], [r_t - r^*] \rangle + 2(c - \frac{1}{c})GR$. The proof is completed. \square

LEMMA A.3. *If the total variation $\max_P \text{Var}_{P(i)}(\theta_i) \leq \sigma^2$, then both the objective function for the zero-sum game L_0 and for the distribution matching game L_{DM} are σ^2 -smooth with respect to r .*

PROOF. (Lemma A.3) Because the Hessian Matrix for the objective functions of both the zero-sum game and the distribution matching game are the same, we use L to represent the objective functions for both games. In other words, the following proof works when we replace L with either L_0 or L_{DM} . To prove σ^2 -smoothness, we need to prove

$$\|\nabla L(r_1) - \nabla L(r_2)\|_2 \leq \sigma^2 \|r_1 - r_2\|_2, \forall r_1, r_2 \in \text{dom } L.$$

Because of the mean value theorem, there exists a point $\tilde{r} \in (r_1, r_2)$ such that

$$\nabla L(r_1) - \nabla L(r_2) = \nabla(\nabla L(\tilde{r}))(r_1 - r_2).$$

Taking the L_2 norm for both sides, we have

$$\begin{aligned} \|\nabla L(r_1) - \nabla L(r_2)\|_2 &= \|\nabla(\nabla L(\tilde{r}))(r_1 - r_2)\|_2 \\ &\leq \|\nabla(\nabla L(\tilde{r}))\|_2 \|r_1 - r_2\|_2 \end{aligned} \quad (33)$$

Then, the problem is to bound the matrix 2-norm $\|\nabla(\nabla L(\tilde{r}))\|_2$. We know in both cases of zero-sum game and distribution matching game:

$$\begin{aligned}\nabla(\nabla L(r)) &= \sum_{i=1}^N P(i)\theta_i\theta_i^T - \left(\sum_{i=1}^N P(i)\theta_i\right)\left(\sum_{j=1}^N P(j)\theta_j\right)^T \\ &= \text{Cov}(\theta_i, \theta_j).\end{aligned}\quad (34)$$

We see $\nabla(\nabla L(r))$ is in the form of a co-variance matrix, which is both symmetric and positive semi-definite. Notice here the probability $P(i)$ depends on r . According to matrix theory, the 2-norm of the matrix is its largest eigenvalue λ_{max} . Further because the covariance matrix is positive semi-definite, all its eigenvalues are non-negative. Hence:

$$\lambda_{max} \leq \sum_i \lambda_i = \text{Tr}(\text{Cov}(\theta_i, \theta_j)).$$

where $\text{Tr}(\text{Cov}(\theta_i, \theta_j))$ means the trace of the covariance matrix (i.e., the sum of all its diagonal entries). Write out the trace, we found it is exactly the total variation. Hence, we have:

$$\begin{aligned}\|\nabla(\nabla L(\tilde{r}))\|_2 &= \lambda_{max} \leq \text{Tr}(\text{Cov}(\theta_i, \theta_j)) \\ &= \mathbb{E}_{P(i)}[\|\theta_i\|_2^2] - \|\mathbb{E}_{P(i)}[\theta_i]\|_2^2 \\ &= \text{Var}_{P(i)}(\theta_i) \leq \sigma^2.\end{aligned}$$

Combining this with Equation 33, we know

$$\|\nabla L(r_1) - \nabla L(r_2)\|_2 \leq \sigma^2 \|r_1 - r_2\|_2.$$

This completes the proof. \square

PROOF. (Theorem 3.1) In Algorithm 2, the gradient g_t is estimated using the mean of K samples $\frac{1}{K} \sum_{i=1}^K \theta_{l'_i}$ according to the gradient in Theorem 2.2. For every sample l'_i sampled from the approximated distribution $P'(l'_i)$, the following inequalities hold because of XOR sampling (Theorem 2.4),

$$\begin{aligned}\frac{1}{\delta} \mathbb{E}_{P(i)}[\theta_i^+] &\leq \mathbb{E}_{P'(l'_i)}[\theta_{l'_i}^+] \leq \delta \mathbb{E}_{P(i)}[\theta_i^+], \\ \delta \mathbb{E}_{P(i)}[\theta_i^-] &\leq \mathbb{E}_{P'(l'_i)}[\theta_{l'_i}^-] \leq \frac{1}{\delta} \mathbb{E}_{P(i)}[\theta_i^-].\end{aligned}$$

To complete the proof, we need to bound the variance of gradient estimation: $\text{Var}(g_t) = \text{Var}(\frac{1}{K} \sum_{i=1}^K \theta_{l'_i})$.

$$\begin{aligned}\|\mathbb{E}_{P(i)}[\theta_i]\|_2^2 &= \|\mathbb{E}_{P(i)}[\theta_i^+ + \theta_i^-]\|_2^2 \\ &= \|\mathbb{E}_{P(i)}[\theta_i^+]\|_2^2 + \|\mathbb{E}_{P(i)}[\theta_i^-]\|_2^2 + \\ &\quad 2\langle \mathbb{E}_{P(i)}[\theta_i^+], \mathbb{E}_{P(i)}[\theta_i^-] \rangle \\ &\leq \|\mathbb{E}_{P(i)}[\theta_i^+]\|_2^2 + \|\mathbb{E}_{P(i)}[\theta_i^-]\|_2^2 \\ &\leq 2G^2\end{aligned}$$

$$\begin{aligned}\text{Var}(\theta_{l'_i}) &= \mathbb{E}_{P'(l'_i)}[\|\theta_{l'_i}\|_2^2] - \|\mathbb{E}_{P'(l'_i)}[\theta_{l'_i}]\|_2^2 \\ &\leq \mathbb{E}_{P'(l'_i)}[\|\theta_{l'_i}\|_2^2] \\ &\leq \delta \mathbb{E}_{P(i)}[\|\theta_i\|_2^2] \\ &= \delta(\text{Var}_{P(i)}(\theta_i) + \|\mathbb{E}_{P(i)}[\theta_i]\|_2^2) \\ &\leq \delta(\sigma^2 + 2G^2)\end{aligned}$$

$$\text{Var}(g_t) = \text{Var}\left(\frac{1}{K} \sum_{i=1}^K \theta_{l'_i}\right) = \frac{1}{K} \text{Var}(\theta_{l'_i}) \leq \frac{\delta}{K}(\sigma^2 + 2G^2)$$

By σ^2 -smooth of L_0 , for the t -th SGD iteration,

$$\begin{aligned}L_0(r_{t+1}) &\leq L_0(r_t) + \langle \nabla L_0(r_t), r_{t+1} - r_t \rangle + \frac{\sigma^2}{2} \|r_{t+1} - r_t\|_2^2 \\ &= L_0(r_t) - \eta \langle \nabla L_0(r_t), g_t \rangle + \frac{\sigma^2 \eta^2}{2} \|g_t\|_2^2\end{aligned}$$

Take the expectation w.r.t. g_t on both sides,

$$\mathbb{E}[L_0(r_{t+1})] \leq L_0(r_t) - \eta \langle \nabla L_0(r_t), \mathbb{E}[g_t] \rangle + \frac{\sigma^2 \eta^2}{2} \mathbb{E}[\|g_t\|_2^2]$$

Before getting the bound of $\|\mathbb{E}[g_t]\|_2^2$, we need to bound $\|\mathbb{E}[g_t^+]\|_2$ and $\|\mathbb{E}[g_t^-]\|_2$ as required in Lemma A.2. Notice that $\mathbb{E}[g_t] = \mathbb{E}[\frac{1}{K} \sum_{i=1}^K \theta_{l'_i}] = \mathbb{E}_{P'(i)}[\theta_i]$. According to Theorem 2.4,

$$\begin{aligned}\mathbb{E}[g_t^+] &= \mathbb{E}_{P'(i)}[\theta_i^+] = \sum_i P'(i)\theta_i^+ \leq \delta \sum_i P(i)\theta_i^+ = \delta \mathbb{E}_{P(i)}[\theta_i^+] \\ \mathbb{E}[g_t^-] &= \mathbb{E}_{P'(i)}[\theta_i^-] = \sum_i P'(i)\theta_i^- \geq \delta \sum_i P(i)\theta_i^- = \delta \mathbb{E}_{P(i)}[\theta_i^-]\end{aligned}$$

Since all entries of $\mathbb{E}[g_t^+]$ are larger or equal to zero, and all entries of $\mathbb{E}[g_t^-]$ are less or equal to zero, calculate the l_2 -norm on both sides of the inequalities above. We have $\|\mathbb{E}[g_t^+]\|_2 \leq \delta \|\mathbb{E}_{P(i)}[\theta_i^+]\|_2 = \delta G$ and $\|\mathbb{E}[g_t^-]\|_2 \leq \delta \|\mathbb{E}_{P(i)}[\theta_i^-]\|_2 = \delta G$. Notice that in Lemma A.2, it was proved that

$$\langle \nabla L_0(r_t), \mathbb{E}[g_t] \rangle \geq \frac{1}{\delta} \|\mathbb{E}[g_t]\|_2^2 - 2(\delta^3 - \delta)G^2$$

Thus we have,

$$\begin{aligned}\mathbb{E}[L_0(r_{t+1})] &\leq L_0(r_t) - \frac{\eta}{\delta} \|\mathbb{E}[g_t]\|_2^2 + \frac{\sigma^2 \eta^2}{2} \mathbb{E}[\|g_t\|_2^2] + 2\eta(\delta^3 - \delta)G^2 \\ &= L_0(r_t) - \frac{\eta}{\delta} (\mathbb{E}[\|g_t\|_2^2] - \text{Var}(g_t)) + \frac{\sigma^2 \eta^2}{2} \mathbb{E}[\|g_t\|_2^2] + 2\eta(\delta^3 - \delta)G^2 \\ &\leq L_0(r_t) - \frac{2\eta - \delta \sigma^2 \eta^2}{2\delta} \mathbb{E}[\|g_t\|_2^2] + \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta(\delta^3 - \delta)G^2 \\ &\leq L_0(r_t) - \frac{\eta \delta}{2} \mathbb{E}[\|g_t\|_2^2] + \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta(\delta^3 - \delta)G^2\end{aligned}$$

Since L_0 is convex w.r.t. r , i.e. $L_0(r^*) \geq L_0(r_t) + \langle \nabla L_0(r_t), r^* - r_t \rangle$, and Lemma A.2 shows $\langle \nabla L_0(r_t), r_t - r^* \rangle \leq \delta \langle \nabla \mathbb{E}[g_t], r_t - r^* \rangle +$

$2(\delta^2 - 1)GR$, we have,

$$\begin{aligned}\mathbb{E}[L_0(r_{t+1})] &\leq L_0(r^*) + \langle \nabla L_0(r_t), r_t - r^* \rangle - \frac{\eta\delta}{2}\mathbb{E}[\|g_t\|_2^2] + \\ &\quad \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta(\delta^3 - \delta)G^2 \\ &\leq L_0(r^*) + \delta\langle \mathbb{E}[g_t], r_t - r^* \rangle - \frac{\eta\delta}{2}\mathbb{E}[\|g_t\|_2^2] + \\ &\quad \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta(\delta^3 - \delta)G^2 + 2(\delta^2 - 1)GR \\ &= L_0(r^*) + \delta\mathbb{E}[\langle g_t, r_t - r^* \rangle - \frac{\eta}{2}\|g_t\|_2^2] + \\ &\quad \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta(\delta^3 - \delta)G^2 + 2(\delta^2 - 1)GR\end{aligned}$$

Then rewrite the term $\delta\mathbb{E}[\langle g_t, r_t - r^* \rangle - \frac{\eta}{2}\|g_t\|_2^2]$ by completing the square. Let

$$\Lambda = \frac{\eta}{K}(\sigma^2 + 2G^2) + 2\eta(\delta^3 - \delta)G^2 + 2(\delta^2 - 1)GR$$

Then we have,

$$\begin{aligned}\mathbb{E}[L_0(r_{t+1})] &\leq L_0(r^*) + \frac{\delta}{2\eta}\mathbb{E}[2\eta\langle g_t, r_t - r^* \rangle - \eta^2\|g_t\|_2^2] + \Lambda \\ &\leq L_0(r^*) + \frac{\delta}{2\eta}\mathbb{E}[\|r_t - r^*\|_2^2 - \|r_t - r^* - \eta g_t\|_2^2] + \Lambda \\ &= L_0(r^*) + \frac{\delta}{2\eta}\mathbb{E}[\|r_t - r^*\|_2^2 - \|r_{t+1} - r^*\|_2^2] + \Lambda\end{aligned}$$

Summing the equations above for $t = 0, \dots, T-1$,

$$\begin{aligned}\sum_{t=0}^{T-1} \mathbb{E}[L_0(r_{t+1}) - L_0(r^*)] &\leq \frac{\delta}{2\eta}\mathbb{E}[\|r_0 - r^*\|_2^2 - \|r_T - r^*\|_2^2] + T\Lambda \\ &\leq \frac{\delta\|r_0 - r^*\|_2^2}{2\eta} + T\Lambda.\end{aligned}$$

According to Jensen's inequality $TL_0(\bar{r}_T) \leq \sum_{t=1}^T L_0(r_t)$, we have

$$\begin{aligned}\sum_{t=0}^{T-1} \mathbb{E}[L_0(r_{t+1}) - L_0(r^*)] &= \mathbb{E}[\sum_{t=1}^T L_0(r_t)] - TL_0(r^*) \\ &\geq T\mathbb{E}[L_0(\bar{r}_T)] - TL_0(r^*)\end{aligned}$$

Combining the equations above, we have

$$\begin{aligned}\mathbb{E}[L_0(\bar{r}_T)] &\leq L_0(r^*) + \frac{\delta\|r_0 - r^*\|_2^2}{2\eta T} + \frac{\eta}{K}(\sigma^2 + 2G^2) + \\ &\quad 2\eta(\delta^3 - \delta)G^2 + 2(\delta^2 - 1)GR.\end{aligned}$$

The proof is completed. \square

A.3 Proofs for Distribution Matching Leader Follower Games

The XOR-Game algorithm for solving distribution matching leader follower game is shown in Algorithm 2. The performance of this algorithm is guaranteed by Theorem 3.3 and Theorem 3.4. To prove Theorem 3.4, we need the following lemmas:

PROOF. (Lemma 3.5) To prove inequality 10, without losing generality, suppose the i -dimension of $\nabla p(r_t)$, namely $[\nabla p(r_t)]_i$ is positive. Because k_t and $\nabla p(r_t)$ match signs at every dimension,

we must have the i -th dimension of k_t , namely $[k_t]_i$ also positive. Thus $\mathbb{E}[k_t]_i$ is also positive. Further because $\frac{1}{c}[\nabla p(r_t)]^+ \leq \mathbb{E}[k_t^+] \leq c[\nabla p(r_t)]^+$, we have inequality 10 holds for dimension i . Suppose the j -dimension of $\nabla p(r_t)$, namely $[\nabla p(r_t)]_j$ is negative. Using similar argument we have $\mathbb{E}[k_t]_j$ is also negative. Now using $c[\nabla p(r_t)]^- \leq \mathbb{E}[k_t^-] \leq \frac{1}{c}[\nabla p(r_t)]^-$, we have inequality 10 holds for dimension j as well. Following similar proofs, one can prove inequalities 11, 12, 13. \square

PROOF. (Lemma 3.6) To prove Equation 14, notice

$$\mathbb{E}[\langle k_t, l_t \rangle] - \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle = \sum_i \mathbb{E}[[k_t]_i [l_t]_i] - \mathbb{E}[[k_t]_i] \mathbb{E}[[l_t]_i]$$

Here, $[k_t]_i([l_t]_i)$ means the i -th dimension of $k_t(l_t)$. Hence, inside the summation of the previous equation is the covariance of $[k_t]_i$ and $[l_t]_i$, namely $\text{cov}([k_t]_i, [l_t]_i)$. We know that the Pearson correlation coefficient ρ_i satisfies:

$$-1 \leq \rho_i = \frac{\text{cov}([k_t]_i, [l_t]_i)}{\sqrt{\text{Var}([k_t]_i)}\sqrt{\text{Var}([l_t]_i)}} \leq 1. \quad (35)$$

Hence:

$$\begin{aligned}&|\mathbb{E}[\langle k_t, l_t \rangle] - \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle| \\ &= \left| \sum_i \text{cov}([k_t]_i, [l_t]_i) \right| \\ &\leq \sum_i |\text{cov}([k_t]_i, [l_t]_i)| \\ &\leq \sum_i \sqrt{\text{Var}([k_t]_i)}\sqrt{\text{Var}([l_t]_i)} \\ &= \sum_i \sqrt{\mathbb{E}[[k_t]_i^2] - (\mathbb{E}[[k_t]_i])^2} \sqrt{\mathbb{E}[[l_t]_i^2] - (\mathbb{E}[[l_t]_i])^2}\end{aligned}$$

Apply the Cauchy-Schwarz Inequality to the last line,

$$\begin{aligned}&|\mathbb{E}[\langle k_t, l_t \rangle] - \langle \mathbb{E}[k_t], \mathbb{E}[l_t] \rangle| \\ &\leq \sqrt{\left(\sum_i \mathbb{E}[[k_t]_i^2] - (\mathbb{E}[[k_t]_i])^2 \right)} \cdot \sqrt{\left(\sum_i \mathbb{E}[[l_t]_i^2] - (\mathbb{E}[[l_t]_i])^2 \right)} \\ &= \sqrt{\text{Var}(k_t)}\sqrt{\text{Var}(l_t)} \\ &\leq \sigma^2.\end{aligned}$$

To prove inequality 15, notice that

$$\begin{aligned}\mathbb{E}[\langle k_t, l_t \rangle] &= \mathbb{E}\left[\sum_i [k_t]_i [l_t]_i\right] \\ &\leq \frac{1}{2}\mathbb{E}\left[\sum_i [k_t]_i^2 + [l_t]_i^2\right] \\ &= \frac{1}{2}(\mathbb{E}[\|k_t\|_2^2] + \mathbb{E}[\|l_t\|_2^2]) \\ &= \frac{1}{2}(\text{Var}(k_t) + \|\mathbb{E}[k_t]\|_2^2 + \text{Var}(l_t) + \|\mathbb{E}[l_t]\|_2^2) \\ &\leq \frac{1}{2}(\sigma^2 + G^2 + \sigma^2 + G^2) = \sigma^2 + G^2.\end{aligned}$$

The first inequality in the chain above is due to the inequality of arithmetic and geometric means. \square

PROOF. (Lemma 3.7) Notice that

$$\langle \nabla f(r_t), r_t - r^* \rangle = \langle \nabla p(r_t) - \nabla q(r_t), r_t - r^* \rangle$$

We split our discussions on multiple conditions concerning the signs of $\nabla p(r_t)$, $\nabla q(r_t)$ and $r_t - r^*$ at each dimension. Notice that $\nabla p(r_t)$, $\nabla q(r_t)$ match signs at each dimension. Hence they are either both positive or negative. In addition, the signs of $r_t - r^*$ at each dimension match that of $\nabla f(r_t) = \nabla p(r_t) - \nabla q(r_t)$ because of the convexity of $f(r_t)$.

First case, suppose at dimension i_1 , $[\nabla p(r_t)]_{i_1}$, $[\nabla q(r_t)]_{i_1}$ and $[r_t - r^*]_{i_1}$ are all positive. In this dimension, under the condition of Theorem 3.4, $[\nabla p(r_t)]_{i_1} \leq c\mathbb{E}[k_t]_{i_1}$ and $[\nabla q(r_t)]_{i_1} \geq \frac{1}{c}\mathbb{E}[l_t]_{i_1}$. Multiply with the positive $[r_t - r^*]_{i_1}$, we have

$$\begin{aligned}
& \langle [\nabla p(r_t) - \nabla q(r_t)]_{i_1}, [r_t - r^*]_{i_1} \rangle \\
& \leq \langle c\mathbb{E}[k_t]_{i_1} - \frac{1}{c}\mathbb{E}[l_t]_{i_1}, [r_t - r^*]_{i_1} \rangle \\
& = c\langle \mathbb{E}[k_t]_{i_1} - \mathbb{E}[l_t]_{i_1}, [r_t - r^*]_{i_1} \rangle + \left(c - \frac{1}{c}\right) \langle \mathbb{E}[l_t]_{i_1}, [r_t - r^*]_{i_1} \rangle \\
& \leq c\langle \mathbb{E}[k_t]_{i_1} - \mathbb{E}[l_t]_{i_1}, [r_t - r^*]_{i_1} \rangle + \\
& \quad \left(c - \frac{1}{c}\right) \langle \max\{|\mathbb{E}[k_t]_{i_1}|, |\mathbb{E}[l_t]_{i_1}|\}, |r_t - r^*|_{i_1} \rangle \quad (36)
\end{aligned}$$

Our second case is when $[\nabla p(r_t)]_{i_2}$, $[\nabla q(r_t)]_{i_2}$ and $[r_t - r^*]_{i_2}$ are all negative. In this case we use $[\nabla p(r_t)]_{i_2} \geq c\mathbb{E}[k_t]_{i_2}$ and $[\nabla q(r_t)]_{i_2} \leq \frac{1}{c}\mathbb{E}[l_t]_{i_2}$ and multiply with the negative term $[r_t - r^*]_{i_2}$, and follow the same derivation as in the first case (except for switching the directions of inequalities when multiplying with negative numbers). We arrive at the same bound as in Equation 36. The bound still holds because we take the absolute values of the last few terms.

The third and the fourth cases are when $[\nabla p(r_t)]_{i_3(i_4)}$, $[\nabla q(r_t)]_{i_3(i_4)}$ are positive (negative) but $[r_t - r^*]_{i_3(i_4)}$ are negative (positive). Following a similar derivation as previous cases,

$$\begin{aligned}
& \langle [\nabla p(r_t) - \nabla q(r_t)]_{i_3(i_4)}, [r_t - r^*]_{i_3(i_4)} \rangle \\
& \leq \langle \frac{1}{c}\mathbb{E}[k_t]_{i_3(i_4)} - c\mathbb{E}[l_t]_{i_3(i_4)}, [r_t - r^*]_{i_3(i_4)} \rangle \\
& = c\langle \mathbb{E}[k_t]_{i_3(i_4)} - \mathbb{E}[l_t]_{i_3(i_4)}, [r_t - r^*]_{i_3(i_4)} \rangle - \\
& \quad \left(c - \frac{1}{c}\right) \langle \mathbb{E}[k_t]_{i_3(i_4)}, [r_t - r^*]_{i_3(i_4)} \rangle \\
& \leq c\langle \mathbb{E}[k_t]_{i_3(i_4)} - \mathbb{E}[l_t]_{i_3(i_4)}, [r_t - r^*]_{i_3(i_4)} \rangle + \\
& \quad \left(c - \frac{1}{c}\right) \langle \max\{|\mathbb{E}[k_t]_{i_3(i_4)}|, |\mathbb{E}[l_t]_{i_3(i_4)}|\}, |r_t - r^*|_{i_3(i_4)} \rangle \quad (37)
\end{aligned}$$

We see that the same bound as the previous two cases can be obtained in the last line. Summing up bounds in Equations 36 and 37 over all dimensions, the left-hand side becomes $\langle \nabla f(r_t), r_t - r^* \rangle$, the first term in the right-hand side becomes $c\langle \mathbb{E}[k_t] - \mathbb{E}[l_t], r_t - r^* \rangle$. The second term in the right-hand side becomes

$\left(c - \frac{1}{c}\right) \sum_i \langle \max\{|\mathbb{E}[k_t]_i|, |\mathbb{E}[l_t]_i|\}, |r_t - r^*|_i \rangle$. Using Cauchy-Schwarz

Inequality for the second term, we get:

$$\begin{aligned}
& \left(c - \frac{1}{c}\right) \sum_i \langle \max\{|\mathbb{E}[k_t]_i|, |\mathbb{E}[l_t]_i|\}, |r_t - r^*|_i \rangle \\
& \leq \left(c - \frac{1}{c}\right) \sqrt{\sum_i \max\{(\mathbb{E}[k_t]_i)^2, (\mathbb{E}[l_t]_i)^2\}} \sqrt{\sum_i [r_t - r^*]_i^2} \\
& \leq \left(c - \frac{1}{c}\right) \sqrt{\sum_i (\mathbb{E}[k_t]_i)^2 + (\mathbb{E}[l_t]_i)^2} \sqrt{\sum_i [r_t - r^*]_i^2} \\
& = \left(c - \frac{1}{c}\right) \sqrt{\|\mathbb{E}[k_t]\|_2^2 + \|\mathbb{E}[l_t]\|_2^2} \sqrt{\sum_i [r_t - r^*]_i^2} \\
& \leq \sqrt{2} \left(c - \frac{1}{c}\right) GR.
\end{aligned}$$

Hence, inequality 16 in Lemma 3.7 holds. \square

PROOF. (Theorem 3.3) In Algorithm XOR-Game, we use the mean of K samples $1/K \sum_{i=1}^K \theta'_i$ to approximate the first part of the gradient in Equation 7: $\mathbb{E}_{P(i)}[\theta_i] = \sum_{i=1}^N P(i)\theta_i$. Here, each l'_i is sampled using XOR sampling, from an approximate probability distribution $P'(l'_i)$. The true distribution is P . According to Theorem 2.4, for any sample l'_i ,

$$\frac{1}{\delta} P(l'_i) \leq P'(l'_i) \leq \delta P(l'_i). \quad (38)$$

Also according to Theorem 2.4, we have:

$$\begin{aligned}
\frac{1}{\delta} \mathbb{E}_{P(i)}[\theta_i^+] & \leq \mathbb{E}_{P'(l'_i)}[\theta_i^+] \leq \delta \mathbb{E}_{P(i)}[\theta_i^+], \\
\delta \mathbb{E}_{P(i)}[\theta_i^-] & \leq \mathbb{E}_{P'(l'_i)}[\theta_i^-] \leq \frac{1}{\delta} \mathbb{E}_{P(i)}[\theta_i^-].
\end{aligned}$$

This implies:

$$\begin{aligned}
\frac{1}{\delta} \mathbb{E}_{P(i)}[\theta_i^+] & \leq \mathbb{E}_{P'(l'_i)} \left[\frac{1}{K} \sum_{i=1}^K \theta_i^+ \right] \leq \delta \mathbb{E}_{P(i)}[\theta_i^+], \\
\delta \mathbb{E}_{P(i)}[\theta_i^-] & \leq \mathbb{E}_{P'(l'_i)} \left[\frac{1}{K} \sum_{i=1}^K \theta_i^- \right] \leq \frac{1}{\delta} \mathbb{E}_{P(i)}[\theta_i^-].
\end{aligned}$$

Similar bounds can be established for $1/S \sum_{j=1}^S \theta''_j$, which is used to approximate the second part of the gradient in Equation 7: $\mathbb{E}_{Q(j)}[\theta_j] = \sum_{j=1}^N Q(j)\theta_j$.

Because $\{\theta_1, \dots, \theta_N\}$ match signs at every dimension, $1/K \sum_{i=1}^K \theta'_i$, $\mathbb{E}_{P(i)}[\theta_i]$, $1/S \sum_{j=1}^S \theta''_j$, and $\mathbb{E}_{Q(j)}[\theta_j]$ all match signs at every dimension. In order to apply Theorem 3.4, we need to bound $\text{Var}(1/K \sum_{i=1}^K \theta'_i)$ and $\text{Var}(1/S \sum_{j=1}^S \theta''_j)$. We have:

$$\begin{aligned}
\text{Var}(l'_i) & = \mathbb{E}_{P'(l'_i)}[\|\theta'_i\|_2^2] - \|\mathbb{E}_{P'(l'_i)}[\theta'_i]\|_2^2 \\
& \leq \mathbb{E}_{P'(l'_i)}[\|\theta'_i\|_2^2] \\
& \leq \delta \mathbb{E}_{P(i)}[\|\theta_i\|_2^2] \\
& = \delta(\text{Var}_{P(i)}(\theta_i) + \|\mathbb{E}_{P(i)}[\theta_i]\|_2^2) \\
& \leq \delta(\sigma^2 + G^2).
\end{aligned}$$

The derivation from the second to the third equation is due to Equation 38. Because $\text{Var}(1/K \sum_{i=1}^K \theta'_i) = 1/K \text{Var}(l'_i)$, we have $\text{Var}(1/K \sum_{i=1}^K \theta'_i) \leq \delta(\sigma^2 + G^2)/K$. Similarly, $\text{Var}(1/S \sum_{j=1}^S \theta''_j) \leq$

$\delta(\sigma^2 + G^2)/S$. Because Equation 38 and all $\{\theta_1, \dots, \theta_N\}$ match signs at every dimension, we have

$$\|\mathbb{E}_{P'(l'_i)}[\theta_{l'_i}]\|_2^2 \leq \delta^2 \|\mathbb{E}_{P(i)}[\theta_i]\|_2^2 = \delta^2 G^2.$$

Because l'_1, \dots, l'_K are i.i.d. sampled, $1/K \sum_{i=1}^K \theta_{l'_i}$ has the same expectation as l'_1 . Hence:

$$\left\| \mathbb{E}_{P'(l'_i)} \left[\frac{1}{K} \sum_{i=1}^K \theta_{l'_i} \right] \right\|_2^2 \leq \delta^2 G^2.$$

Similarly, $\|\mathbb{E}_{Q'(l''_j)}(\frac{1}{S} \sum_{j=1}^S \theta_{l''_j})\|_2^2 \leq \delta^2 G^2$. Apply all the bounds computed above into Equation 9 in Theorem 3.4, also notice L is σ^2 -smooth due to Lemma A.3, we get the following bound:

$$\begin{aligned} \mathbb{E}[L(\overline{r_T})] - L(r^*) &\leq \frac{\delta \|r_0 - r^*\|_2^2}{2\eta T} + \\ (\delta^2 - 1) \left[\sqrt{2GR} + 2\eta \left(\frac{\sigma^2 + G^2}{\min\{K, S\}} + \delta G^2 \right) \right] &+ 2\eta(\delta^2 + 1) \frac{\sigma^2 + G^2}{\min\{K, S\}}. \end{aligned}$$

Proof complete. \square

A.4 Proofs for Number of NP Oracles Needed

The proof for the number of NP oracles needed is developed from [16]. We encourage the readers to read the original papers for better understanding.

PROOF. (Theorem 3.8) From Theorem 2.4, it requires $O(-n \log(1 - 1/\sqrt{\delta}) \log \frac{-n \log(1 - 1/\sqrt{\delta})}{\gamma})$ queries of NP oracles to generate one sample. However, as specified in [16], once we have the first sample, the following samples will not need as many queries. Therefore, generating K samples can be seen as generating one sample first, then generating following samples inside the same SGD iteration next. We fix the number of XOR constraints needed to be added starting the generation of the second sample (in other words, the ComputeK procedure in [16] can be avoided for the rest $K - 1$ samples). As a result, we need $O(K - 1)$ NP oracle queries in generating each of the following $K - 1$ samples. Thus the total queries for K samples will be $O(-n \log(1 - 1/\sqrt{\delta}) \log \frac{-n \log(1 - 1/\sqrt{\delta})}{\gamma} + K)$. To complete T SGD iterations, XOR-Game₀ requires $O(-Tn \log(1 - 1/\sqrt{\delta}) \log \frac{-n \log(1 - 1/\sqrt{\delta})}{\gamma} + TK)$ queries to NP oracles. XOR-Game_{DM} needs additional samples from Q , hence overall it needs $O(-n \log(1 - 1/\sqrt{\delta}) \log \frac{-n \log(1 - 1/\sqrt{\delta})}{\gamma} + TK + S)$ queries to NP oracles. \square

B EXPERIMENTAL DETAILS

Here we show additional details for the experiments in the main text. In all experiments, XOR-Game is implemented with CPLEX 12.6. All experiments are run on computational nodes with dual 64-core AMD Epyc 7662 ‘‘Rome’’ processors@2.0GHz with a maximum 256 GB of memory.

Evaluation The zero-sum objective L_0 is evaluated by the exact model counter Ace [10]. For distribution matching games, the performance is evaluated by the utility function of the leader, which is

the KL-Divergence between $Q(L)$ and $P(L)$. Notice

$$\begin{aligned} KL(Q||P) &= \sum_{L \in \mathcal{L}} Q(L) \log \left(\frac{Q(L)}{P(L)} \right) \\ &= \sum_{L \in \mathcal{L}} Q(L) (L^T W_q L - w_r r^T L - w_f F L - L^T W L) + \log Z_p - \log Z_q \\ &= \mathbb{E}_Q [L^T W_q L - w_r r^T L - w_f F L - L^T W L] + \log Z_p - \log Z_q \end{aligned}$$

where Z_p and Z_q are partition functions of $P(L)$ and $Q(L)$. $\mathbb{E}_Q [L^T W_q L - w_r r^T L - w_f F L - L^T W L]$ are approximated with the average of 200 samples using XOR-Sampling. Considering XOR-Sampling has a constant approximation bound, we believe the estimation to the expectation is accurate. For synthetic experiment, we used again exact counter Ace to calculate log partition functions $\log Z_p$ and $\log Z_q$. For Avicaching game evaluation, the scale of the problem is beyond exact approaches. We used the a winning solver HAK [25] solver implemented in libDAI [36] to compute the partition functions.

B.1 Synthetic Benchmarks

Our synthetic dataset consists of 30 locations. To generate $P(L)$, w_r is set to 0.2. w_f is a 1-by-5 weight vector, each entry of which is randomly drawn from a uniform distribution $U[-0.1, 0.1]$. F is a 5-by-30 matrix, each entry of which is generated from uniform distribution $U[-0.1, 0.1]$. W is a 30-by-30 matrix capturing the inter-dependency of locations in affecting participants’ interests. We intentionally add a few large entries in W to increase the difficulty of benchmarks. In particular, each entry of W , W_{ij} is generated by $z_{ij}(b_{ij} + g_{ij})$, where z_{ij} is uniformly sampled from a Binomial distribution with $p = 0.1$ to serve as a mask, b_{ij} is from uniform distribution $U[-3, 3]$. $g_{ij} = 20(h_{ij} - 0.5)$, where h_{ij} is drawn from a Binomial distribution with $p = 0.5$. The initial reward r_0 is sampled from uniform distribution $U(0.1, 1)$. For distribution matching game, $Q(L)$ is generated similarly as $P(L)$, except for keeping all reward r zero. SGD step size $\eta = 0.1$ are used in both experiments.

BRQR vs. XOR-Game For the comparison between BRQR and XOR-Game on zero-sum games, the leader’s expected utilities during as the optimization progress are shown in Fig 2. BRQR slightly converges faster than XOR-Game. Both approaches converge to close-to-optimal solutions. When N grows larger than 2^{16} , which is common in real-world benchmarks, the BRQR soon become infeasible because it needs to go over all the actions in each iteration.

B.2 Avicaching Game

This section provides the additional experiments we have done for the Avicaching Game and additional details for the experiment setup discussed in the main text. The rough idea behind is (1) learning the behavior model from real-world dataset, and (2) applying XOR-Game_{DM} and other baselines to find the best reward that can alter the behavior model towards a desired target distribution. Finally, we evaluate the performance using the KL-Divergence between the final behavior model and the target model.

B.2.1 Learning Behavior Model. This section describes the experiments we have done for learning the behavior model from Avicaching dataset.

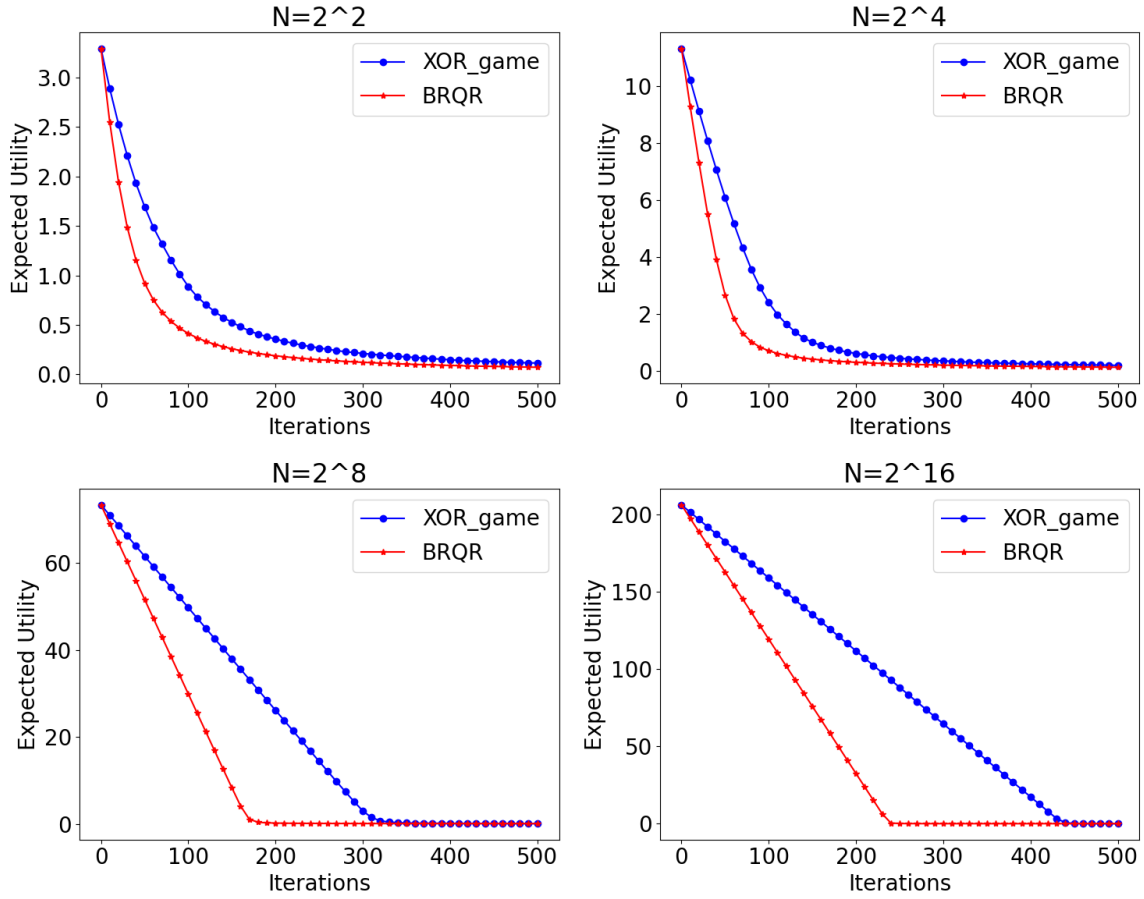


Figure 2: The leader’s expected utility vs. # iterations in different sized games for XOR-Game and BRQR. For small-sized games, BRQR converges in slightly fewer iterations. Both algorithms find close-to-optimal solutions in the end. Combining with Table 1, the actual running time of BRQR grows much faster than XOR-Game as the game sizes increase, and hence cannot scale to games with large action spaces.

Behavior Model As discussed in the main text, the behavior model used in Avicaching game is as follows,

$$P(L|r) = \frac{\exp\left(w_r r^T L + w_f^T F L + L^T W L\right)}{\sum_{L' \in \mathcal{L}} \exp\left(w_r r^T L' + w_f^T F L' + L'^T W L'\right)}$$

$L = [l_1, \dots, l_N]^T \in \{0, 1\}^N$ represents a set of visited locations by setting $l_i = 1$ iff. l_i is visited. $r = [r_1, \dots, r_N]^T \in \mathbb{R}^N$ where r_i represents the reward applied to location i . $w_r \in \mathbb{R}$ is the weighting parameter characterizing the preference of rewards. $w_f \in \mathbb{R}^M$ is the weighting parameter characterizing the preference of the natural land features. $F = [F_1, \dots, F_N] \in \mathbb{R}^{M \times N}$ is the feature matrix containing the land characteristics of all N locations. $W \in \mathbb{R}^{N \times N}$ is the regularization term.

Dataset The dataset for learning the behavior model is generated from the Avicaching game [58] experiment conducted from 03/28/2015 to 10/30/2015 (30 weeks). Participants are encouraged to visit several locations and report their bird observations. Around

50 users participated in this field experiment. There are 116 observation points, noted as \mathcal{L} , of which 50 locations, noted as \mathcal{L}_a , are Avicaching locations. In the actual field experiment, participants are encouraged to visit Avicaching locations by setting rewards ranging from 1.0 to 15.0 while the reward of visiting other locations remain to be 0. The amount of reward changes weekly resulting in 16 different reward schemes.

To generate enough datapoints for training, all participants are seen to be identical. For each participant, the locations visited in one week is seen as his/her set of visited locations under the reward that week. Therefore, we can generate one location set vector L together with the corresponding reward vector r as one data point. We successfully generated 1112 pairs (L, r) .

The feature matrix F contains land characteristics of all observation points from the National Land Cover Database (NLCD) 2011. The database provides spatially explicit and reliable information on the Nation’s land cover within a 375 meter circle around each location. We finally selected 32 features including latitude, longitude, and 30 distinguishable land cover features. The 32-dimensional

feature vectors of 116 locations formulate the feature matrix $F \in \mathbb{R}^{32 \times 116}$.

Learning Setup Suppose we have K data points $\{(L_i, r_i)\}$ with a known feature matrix F , the log likelihood function $l(w_r, w_f, W) = \frac{1}{K} \sum_{i=1}^K \log P(L_i|r)$ is:

$$l(w_r, w_f, W) = \frac{1}{K} \sum_{i=1}^K w_r r^T L_i - \Lambda(w_r, w_f, W)$$

$$\Lambda(w_r, w_f, W) = \log \sum_{L \in \mathcal{L}} e^{w_r r^T L + w_f^T F L + L^T W L}$$

From Equation 6. in XOR-CD. The gradient of $l(w_r, w_f, W)$ with respect to w_r , w_f , and W can be estimated as

$$g_{cd}(w_r) = \frac{1}{K} \sum_{i=1}^K r^T L_i - \frac{1}{S} \sum_{i=1}^S r^T L'_i$$

$$g_{cd}(w_f) = \frac{1}{K} \sum_{i=1}^K F L_i - \frac{1}{S} \sum_{i=1}^S F L'_i$$

$$g_{cd}(W) = \frac{1}{K} \sum_{i=1}^K L_i L_i^T - \frac{1}{S} \sum_{i=1}^S L'_i L'_i^T$$

where $\{L'_1, \dots, L'_S\}$ are the samples from the current model distribution. We use XOR-Sampling to obtain these samples.

With the gradient estimation, we apply SGD to learn the model parameters. The dataset $\{(L, r)\}$ is divided into 16 batches. In each batch, the reward vectors are the same. The number of all data points is 1112, and the number of data in each batch ranges from 33 to 136. The learning rate is fixed as 0.01 and the total number of SGD iterations is 250. For initialization, w_r is uniformly sampled from $U(0, 1)$. w_f is sampled from $U(0.1, 1)$. W is initialized as a symmetric matrix with entries sampled from $U(0.1, 1)$. For XOR-Sampling, we utilized the work from Ermon et al. [16]. and used the same parameter. The number of samples is 80.

Time Consumption In the early stage, each iteration takes around 480 seconds. From iteration 50, XOR-Sampling takes much more time for queries to NP oracles. Each iteration takes up to 15000 seconds.

Evaluation We used principal component analysis (PCA) to empirically visualize the behavior model learned. Figure 3 shows the 2-dimensional plot of the data points in the real-world experiments and the samples obtained from the learned model (using XOR sampling). Both the data points and the samples are visualized using a 2-dimensional PCA in Figure 3. The plot separates on individual reward levels. We can see that the samples obtained from the learned behavior model replicates the original data distribution well.

B.2.2 XOR-Game & Baseline Settings. Here is additional information regarding the experimental details of XOR-Game and other baselines.

Target Model The idea in Avicaching game is to promote the probability of sample-needed locations. We empirically designed the target model as $Q(L)$, where

$$Q(L) = \frac{e^{L^T W_Q L}}{\sum_{L' \in \mathcal{L}} e^{L'^T W_Q L'}}$$

The entry $[W]_{i,j}$ in i -th row, j -th column satisfies:

$$[W]_{i,j} \sim \begin{cases} U(0, 0.05) & i, j \in \mathcal{L}_a, i = j \\ U(0.01, 0.02) & i, j \in \mathcal{L}_a, i \neq j \\ U(-0.01, 0) & i \in \mathcal{L}_a, j \notin \mathcal{L}_a \\ U(-0.05, 0) & i, j \notin \mathcal{L}_a, i = j \\ U(-0.02, -0.01) & i, j \notin \mathcal{L}_a, i \neq j \end{cases}$$

Generally, the probability of visiting Avicaching locations \mathcal{L}_a is much higher than that of other locations in this Q model.

Experimental details The SGD step size is fixed at 0.1. The maximum iterations $T = 100$ which is enough for convergence. The initialization of reward vector is sampled from a uniform distribution $U(5, 10)$. The number of samples from $Q(L)$ is 100, and the number of samples from $P(L|r_t)$ in each iteration is 50 considering the trade-off between running time and accuracy. For XOR-Sampling, parameters are set the same as in [16]. Apart from the sampling method, all baselines share the same settings as XOR-Game. For Gibbs-Game, we use Gibbs sampling after taking 5800 MCMC steps in replace of XOR-Sampling. In BP-Game, the samples are generated from Belief Propagation [42].

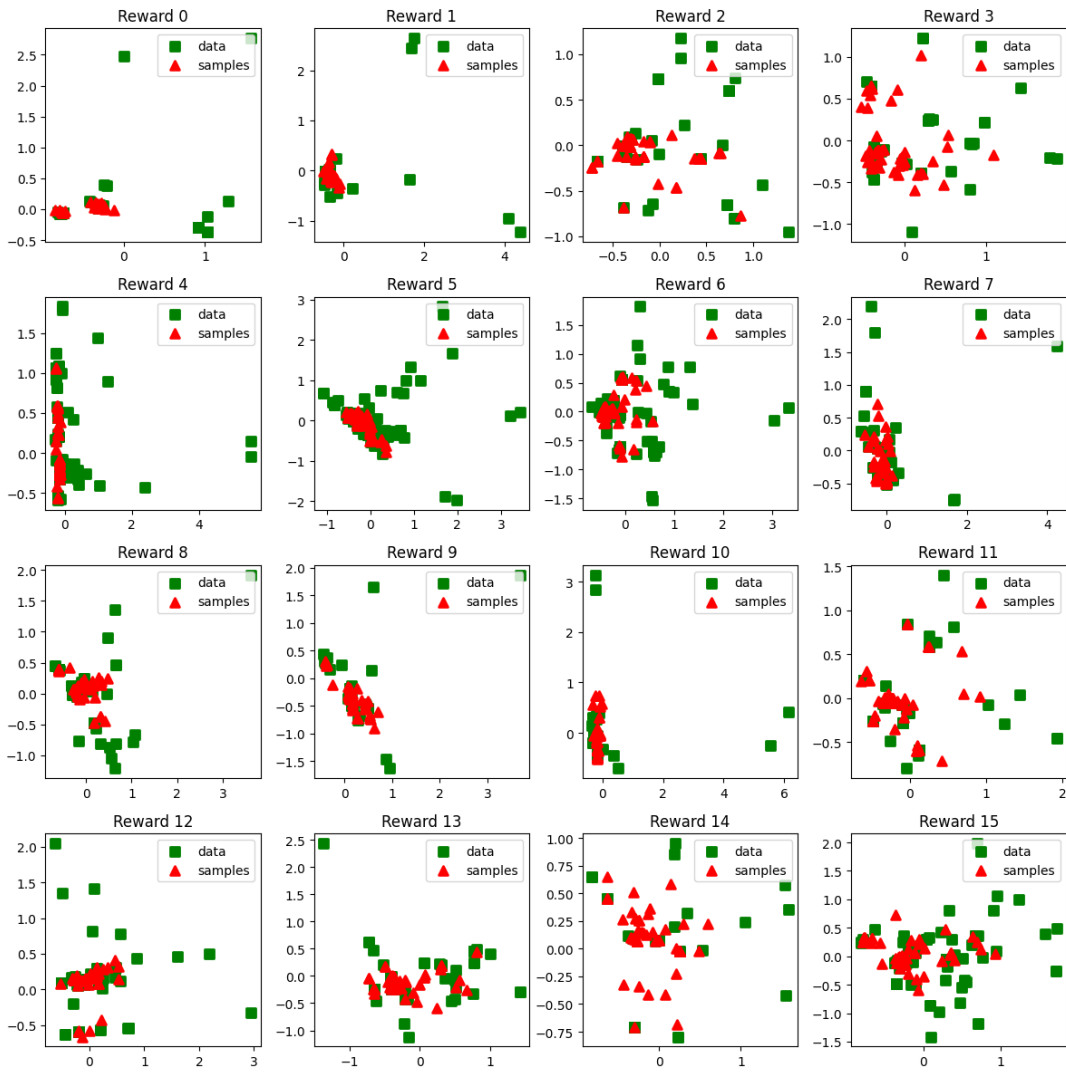


Figure 3: Location set vectors in Avicaching dataset and samples from learned behavior model under different reward level (visualized through 2-dimensional PCA). We can see that the samples obtained from the learned behavior model (red dots) replicates the original data distribution (green dots) well.