

Improved Algorithms for Learning Equilibria in Simulation-Based Games

Enrique Areyan Viqueira
Brown University
Providence, Rhode Island
eareyan@brown.edu

Cyrus Cousins
Brown University
Providence, Rhode Island
cyrus_cousins@brown.edu

Amy Greenwald
Brown University
Providence, Rhode Island
amy_greenwald@brown.edu

ABSTRACT

We tackle a fundamental problem in empirical game-theoretic analysis (EGTA), that of learning equilibria of simulation-based games. Such games cannot be described in analytical form; instead, a black-box simulator can be queried to obtain noisy samples of utilities. Our first theorem establishes that uniform approximations of simulation-based games are equilibrium preserving. We then design algorithms that uniformly approximate simulation-based games with finite-sample guarantees. Our first algorithm, global sampling (GS), extends previous work that constructs confidence intervals assuming bounded utilities with confidence intervals that are sensitive to variance. The second, progressive sample with pruning (PSP), samples progressively, ceasing the sampling process (i.e., pruning strategies) as soon as it determines that the corresponding utilities have been sufficiently well estimated for equilibrium computation. We experiment with our algorithms using both GAMUT, a state-of-the-art game generator, and Gambit, a state-of-the-art game solver. For a broad swath of games, we show that GS using our variance-sensitive bounds outperforms previous work, and that PSP can significantly outperform GS. Here “outperform” means achieving the same guarantees with far fewer samples.

KEYWORDS

Empirical Game-Theoretical Analysis; PAC Learning

ACM Reference Format:

Enrique Areyan Viqueira, Cyrus Cousins, and Amy Greenwald. 2020. Improved Algorithms for Learning Equilibria in Simulation-Based Games. In *Proc. of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2020)*, Auckland, New Zealand, May 9–13, 2020, IFAAMAS, 14 pages.

1 INTRODUCTION

Game theory is the *de facto* standard conceptual framework used to analyze the strategic interactions among rational agents in multi-agent systems. At the heart of this framework is the theoretical notion of a game. In a game, each player (also referred to as an agent) chooses a strategy and earns a utility that in general depends on the profile (i.e., vector) of strategies chosen by all the agents.

The most basic representation of a game is the *normal form*, which can be visualized as a matrix containing all agents’ utilities at all strategy profiles. Analyzing a game means predicting its outcome, i.e., the strategy profiles that one can reasonably expect the agents to play. Perhaps the most common such prediction is

that of *Nash equilibrium*, where each agent plays a *best-response* to the strategies of the others: i.e., a strategy that maximizes their utility [20]. More generally, at an ϵ -Nash equilibrium agents best respond to other agents’ strategies up to an additive error of ϵ .

This paper is concerned with analyzing games for which a complete and accurate description is not available. While knowledge of the number of agents and their respective strategy sets is available, we do not assume *a priori* access to the game’s utility function. Instead, we assume access to a simulator from which we can sample noisy utilities associated with any strategy profile. Such games have been called *simulation-based games* [33] and *black-box games* [22], and their analysis is called *empirical game theoretic analysis* (EGTA) [13, 34]. EGTA methodology has been applied in a variety of practical settings for which simulators are readily available, including trading agent analyses in supply chains [13, 32], ad auctions [1, 15], and energy markets [16]; designing network routing protocols [35]; strategy selection in real-time games [27]; and the dynamics of RL algorithms, like AlphaGo [28].

The aim of this work is to design learning algorithms that can accurately estimate all the equilibria of simulation-based games. We tackle this problem using the *probably approximately correct* (PAC) learning framework [29]. Our algorithms learn so-called *empirical games* [34], which are estimates of simulation-based games constructed via sampling. We argue that empirical games so constructed yield *uniform approximations* of simulation-based games, meaning all utilities in the empirical game tend toward their expected counterparts, *simultaneously*.

This notion of uniform approximation is central to our work: we prove that, when one game Γ' is a uniform approximation of another Γ'' , all equilibria in Γ' are approximate equilibria in Γ'' . In particular, when an estimated game $\hat{\Gamma}$ uniformly approximates a (true) game Γ , all equilibria in Γ are approximate equilibria in $\hat{\Gamma}$, and all equilibria in $\hat{\Gamma}$ are approximate equilibria in Γ . As a result, a uniform approximation implies perfect recall¹ (all true positives—equilibria in Γ —are at least approximate equilibria in $\hat{\Gamma}$) and *approximately* perfect precision (all false positives—equilibria in $\hat{\Gamma}$ but not necessarily in Γ —are at least approximate equilibria in Γ). Our learning algorithms, which learn empirical games that are uniform approximations of simulation-based games, thus well estimate the equilibria of simulation-based games.

Estimating *all* the utilities in an empirical game is non-trivial, in part because of the *multiple comparisons problem* (MCP), which arises when estimating many parameters simultaneously, since accurate inferences for each individual parameter do not necessarily imply a similar degree of accuracy for the parameters in aggregate.

Proc. of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2020), B. An, N. Yorke-Smith, A. El Fallah Seghrouchni, G. Sukthankar (eds.), May 9–13, 2020, Auckland, New Zealand. © 2020 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

¹We use the term *recall* in the information retrieval sense. The juxtaposition of the two words “perfect” and “recall” is not a reference to extensive-form games.

Controlling the *family-wise error rate* (FWER)—the probability that one or more of the approximation guarantees is violated—is one way to control for multiple comparisons. We control the FWER in empirical games by first using a concentration inequality to establish confidence intervals about each parameter *individually*, and then applying a statistical correction (e.g., Bonferroni; in other words, a union bound) to bound the probability that all approximation guarantees hold *simultaneously*.

The first concentration inequality we apply is Hoeffding’s [1963]. We note that Hoeffding’s inequality implies sub-Gaussian tail bounds for averages of m c -bounded random variables, yielding $1 - \delta$ probability confidence intervals of width $\Theta(\sqrt{c^2 \ln(1/\delta)/m})$. Hoeffding’s inequality, however, is sensitive only to the range c of the random variables, not their variance σ^2 . Bennett’s inequality [1962] relaxes the dependency on c , replacing it with a dependency on $\sigma^2 < c^2$, yielding width $\Theta\left(\frac{c \ln(1/\delta)}{m} + \sqrt{\frac{\sigma^2 \ln(1/\delta)}{m}}\right)$; i.e., a sub-gamma tail bound. Tighter still would be to use Gaussian tail bounds, which yield $\Theta(\sqrt{\sigma^2 \ln(1/\delta)/m})$. However, the central limit theorem does not provide finite sample guarantees, and any approximation errors are amplified in a MCP setting; when applied to games, a single mistake can lead to a cascade of mistakes in equilibrium estimation.

Thus, we resort to sub-gamma tail bounds, which still drastically improve the situation because of the presence of a fast-decaying, hyperbolic $O(1/m)$ term in addition to the usual root-hyperbolic $O(1/\sqrt{m})$ term. The fast-decaying range term is initially the dominant factor, but it is quickly overcome by the variance term. When the variance is substantially less than c^2 , these variance-sensitive bounds are substantially tighter than Hoeffding’s. They are also of particular interest when only a loose bound on c is available, or in cases where extreme outliers are unlikely albeit possible, thus avoiding the pitfalls of the central limit theorem. In both cases, c rigorously controls for the possibility of a rare event.

As the variance is not generally known, we present a 2-step strategy wherein the variance is upper-bounded in terms of the empirical variance, and the expectation is then upper-bounded by Bennett’s inequality in terms of this upper bound, appropriately corrected for the MCP. Analogous to the use of Hoeffding’s bound in [28], our new bound (on the expectation) immediately gives rise to a novel algorithm by which to learn an empirical game from a static sample of utilities. We refer to this algorithm generically as *global sampling* (GS), and instantiate with both bounds.

We also describe a second learning algorithm, *progressive sampling with pruning* (PSP), which samples dynamically. Prior art [9, 24, 25] uses progressive sampling to obtain a desired accuracy given a failure probability by guessing an initial sample size, computing Rademacher bounds, and repeating with larger sample sizes until said accuracy is attained. Our work complements this approach with *pruning*: at each iteration, parameters that have already been sufficiently well estimated for the task at hand (here, equilibrium estimation) are pruned as subsequent iterations of PSP do not refine their bounds. Pruning represents both a statistical and computational improvement over earlier progressive sampling techniques.

Finally, we report on extensive experiments designed to evaluate the performance of our algorithms. Producing a robust empirical evaluation is a daunting task, as the space of possible games is vast. Moreover, computing equilibria is computationally intractable [7].

Fortunately, researchers have devised tools to address these issues. The first, GAMUT, is a state-of-the-art suite of game generators capable of producing a myriad of games with rich strategic structure [21]. The second, Gambit, is a state-of-the-art solver to compute Nash equilibria [19]. We use both GAMUT and Gambit to evaluate the performance of our algorithms. Concretely, for a wide variety of game types, we draw multiple random game instances from GAMUT, and for each such instance, we solve for its equilibria using Gambit. Using our variance-sensitive bounds within GS outperforms Hoeffding’s; and PSP, with either bound, can significantly outperform the corresponding GS algorithm. Here “outperform” means achieving the same guarantees with far fewer samples—exponentially fewer for PSP in sufficiently large games (in our experiments, two players with ten or more strategies each).

This paper purports to contribute to the literature on EGTA, a methodology for the analysis of multi-agent systems. One important application of our methodology is the estimation of equilibria in meta-games [27]. *Meta-games* are simplified versions of intractably large games, where, instead of modeling every possible strategy an agent might implement, one analyzes a game with a substantially reduced set of strategies, each of which may be given by a complicated algorithmic procedure (i.e., a heuristic). That is, instead of analyzing a game in terms of its (low-level) formal game-theoretic strategies—a computationally intractable task in a game such as Go—one might analyze a reduced version of the game where agents play according to higher-level strategies given by reinforcement learning algorithms: e.g., variants of AlphaGo [26]. Simulation is in order; and since each run of the game can result in either agent winning (depending on various stochastic elements, including perhaps the agents’ strategies), utilities are noisy in general. Our techniques are directly applicable to the learning of empirical meta-games, and provide finite-sample guarantees on the quality of the equilibria of the corresponding simulation-based meta-games.

Related Work. One distinctive feature of our work *vis à vis* the literature is that we aim to estimate all equilibria of a simulation-based game, rather than just one (e.g., [14]). Notable exceptions include [28], [31], [36], and [2, 30], the latter being our own prior work. The first of these papers argues that all equilibria of a simulation-based game are also approximate equilibria in an empirical game: i.e., they establish perfect recall with finite-sample guarantees. The second paper derives asymptotic results about the quality of the equilibria of empirical games: i.e., they establish perfect precision in the limit, as the number of samples tends to infinity. Our first theorem unifies these two results, obtaining finite sample guarantees for both perfect recall and approximate precision. The third paper is perhaps closest in spirit to ours; the goals are to characterize the quality of the equilibria in empirical games, and to exploit statistical information to save on sampling. Their methods employ bootstrapping [8], so do not immediately afford any theoretical guarantees.

2 APPROXIMATION FRAMEWORK

We begin by presenting standard game-theoretic notions. We then introduce the notion of uniform approximation. Given an approximation of one game by another, there is not necessarily a connection between their properties. For example, there may be equilibria in one game with no corresponding equilibria in the other, as

small changes to the utility functions can add or remove equilibria. Nonetheless, we show that finding the approximate equilibria of a uniform approximation of a game is sufficient for finding all the (exact) equilibria of the game itself.

A *normal-form game* $\Gamma \doteq \langle P, \{S_p \mid p \in P\}, \mathbf{u}(\cdot) \rangle$ consists of a set of agents P , with *pure strategy set* S_p available to agent $p \in P$. We define $S \doteq S_1 \times \cdots \times S_{|P|}$ to be the pure strategy profile space of Γ , and then $\mathbf{u} : S \rightarrow \mathbb{R}^{|P|}$ is a vector-valued utility function.

Given a normal-form game Γ , we let S_p° denote the set of distributions over S_p ; this set contains p 's *mixed strategies*. We define $S^\circ = S_1^\circ \times \cdots \times S_{|P|}^\circ$, and then, overloading notation, we write $\mathbf{u}(s)$ to denote the expected utility of a mixed strategy profile $s \in S^\circ$.

A solution to a normal-form game is a prediction of how strategic agents will play the game. One solution concept that has received a great deal of attention in the literature is Nash equilibrium [20], a (pure or mixed) strategy profile at which each agent selects a utility-maximizing strategy, fixing all other agents' strategies. In this paper, we are concerned with ε -Nash equilibrium, an approximation of Nash equilibrium that is amenable to statistical estimation.

Given a game Γ , fix an agent p and a mixed strategy profile $s \in S^\circ$. We define $T_{p,s}^\circ \doteq \{t \in S^\circ \mid t_q = s_q, \forall q \neq p\}$: i.e., the set of all mixed strategy profiles in which the strategies of all agents $q \neq p$ are fixed at s_q . Agent p 's *regret* at s is defined as $\text{Reg}_p^\circ(\Gamma, s) \doteq \sup_{s' \in T_{p,s}^\circ} \mathbf{u}_p(s') - \mathbf{u}_p(s)$. By restricting s and $T_{p,s}^\circ$ to pure strategy profiles, agent p 's *pure regret* $\text{Reg}_p(\Gamma, s)$ can be defined similarly.

Note that $\text{Reg}_p^\circ(\Gamma, s), \text{Reg}_p(\Gamma, s) \geq 0$, since agent p can deviate to any strategy $s' \in S_p$, including s_p itself. A strategy profile s that has regret at most $\varepsilon \geq 0$, for all $p \in P$, is an ε -Nash equilibrium:

Given $\varepsilon \geq 0$, a mixed strategy profile $s \in S^\circ$ in a game Γ is an ε -Nash equilibrium if, for all $p \in P$, $\text{Reg}_p^\circ(\Gamma, s) \leq \varepsilon$. At a pure strategy ε -Nash equilibrium $s \in S$, for all $p \in P$, $\text{Reg}_p(\Gamma, s) \leq \varepsilon$. Let $E_\varepsilon^\circ(\Gamma)$ ($E_\varepsilon(\Gamma)$) be the set of mixed (pure) ε -Nash equilibria. $E_\varepsilon(\Gamma) \subseteq E_\varepsilon^\circ(\Gamma)$.

We now show that equilibria can be approximated with bounded error, given a uniform approximation. Our main theorem establishes perfect recall: the approximate game contains all true positives: i.e., all (exact) equilibria of the original game. It also establishes approximately perfect precision: all false positives in the approximate game are approximate equilibria in the original game.

We define the ℓ_∞ -norm between two compatible games, with the same agents sets P and strategy profile spaces S , and with utility functions \mathbf{u}, \mathbf{u}' , respectively, as follows: $\|\Gamma - \Gamma'\|_\infty \doteq \|\mathbf{u}(\cdot) - \mathbf{u}'(\cdot)\|_\infty \doteq \sup_{p \in P, s \in S} |\mathbf{u}_p(s) - \mathbf{u}'_p(s)|$. While the ℓ_∞ -norm as defined applies only to pure normal-form games, it is in fact sufficient to use this metric even to show that the utilities of mixed strategy profiles approximate one another. We formalize this claim presently.²

LEMMA 2.1. *If Γ, Γ' differ only in their utilities, then*
 $\sup_{p \in P, s \in S^\circ} |\mathbf{u}_p(s) - \mathbf{u}'_p(s)| = \|\Gamma - \Gamma'\|_\infty$.

Γ' is said to be a *uniform ε -approximation* of Γ when $\|\Gamma - \Gamma'\|_\infty \leq \varepsilon$. They are so-called because the bound between utility deviations in Γ and Γ' holds *uniformly* over all players and strategy profiles.

THEOREM 2.2 (APPROXIMATE EQUILIBRIA). *If two normal-form games, Γ and Γ' , are uniform approximations of one another, then: (1) $E(\Gamma) \subseteq E_{2\varepsilon}(\Gamma') \subseteq E_{4\varepsilon}(\Gamma)$, and (2) $E^\circ(\Gamma) \subseteq E_{2\varepsilon}^\circ(\Gamma') \subseteq E_{4\varepsilon}^\circ(\Gamma)$.*

²All proofs appear in the supplemental material, <http://github.com/eareyan/pysegta>.

3 LEARNING FRAMEWORK

In this section, we move on from approximating equilibria in games to learning them. We present algorithms that learn so-called empirical games, which comprise estimates of the expected utilities of simulation-based games. We further derive uniform convergence bounds, proving that our algorithms output empirical games that uniformly approximate their expected counterparts, with high probability. By Theorem 2.2, the equilibria of these empirical games thus approximate those of the corresponding simulation-based games, with high probability.

A *conditional normal-form game* $\Gamma_X \doteq \langle X, P, \{S_p \mid p \in P\}, \mathbf{u}(\cdot) \rangle$ consists of a set of conditions X , a set of agents P , with pure strategy set S_p available to agent p , and a vector-valued conditional utility function $\mathbf{u} : S \times X \rightarrow \mathbb{R}^{|P|}$. It is convenient to imagine a condition $x \in X$ as pertaining to the set $P \times S$. Given such a condition, $\mathbf{u}(\cdot; x)$ yields a standard utility function of the form $S \rightarrow \mathbb{R}^{|P|}$. Given a conditional normal-form game Γ_X together with distribution \mathcal{D} , we also define the *expected utility function* $\mathbf{u}(s; \mathcal{D}) = \mathbb{E}_{x \sim \mathcal{D}} [\mathbf{u}(s; x)]$, and the *expected normal-form game* as $\Gamma_\mathcal{D} \doteq \langle P, \{S_p \mid p \in P\}, \mathbf{u}(\cdot; \mathcal{D}) \rangle$.

Expected normal-form games serve as our mathematical model of simulation-based games. They are sufficient not only to model arbitrary black-box games, but additionally games where the *rules* are known but *environmental conditions* are random (e.g., auctions where bidder valuations are random, or *deterministic* war games where initial armies and terrain are random), as well as games with *randomness*, where X is taken to be a *PRNG seed* or *entropy source*.

Given a conditional normal-form game Γ_X together with a distribution \mathcal{D} and sample conditions $X = (x_1, \dots, x_m) \sim \mathcal{D}^m$, we define the *empirical utility function* $\hat{\mathbf{u}}(s; X) \doteq \frac{1}{m} \sum_{j=1}^m \mathbf{u}(s; x_j)$. The corresponding *empirical normal-form game* is then $\hat{\Gamma}_X \doteq \langle P, \{S_p \mid p \in P\}, \hat{\mathbf{u}}(\cdot; X) \rangle$.

Our present goal, then, is to “uniformly learn” empirical games (i.e., obtain uniform convergence guarantees) from finitely many samples. This learning problem is non-trivial because it involves multiple comparisons. Tuyls et al. [28] use Hoeffding’s inequality to estimate a single utility value, and then apply a Sidák correction to estimate all utility values simultaneously, assuming independence among agents’ utilities. Similarly, one can apply a Bonferroni correction (i.e., a union bound) to Hoeffding’s inequality, which does not require independence, but yields a slightly looser bound.

THEOREM 3.1 (FINITE-SAMPLE BOUNDS VIA Hoeffding’s INEQUALITY). *Consider finite, conditional normal-form game Γ_X together with distribution \mathcal{D} and index set $I \subseteq P \times S$ such that for all $x \in X$ and $(p, s) \in I$, it holds that $\mathbf{u}_p(s; x) \in [-c/2, c/2]$, where $c \in \mathbb{R}$. Then, with probability at least $1 - \delta$ over $X \sim \mathcal{D}^m$, it holds that*

$$\sup_{(p,s) \in I} |\mathbf{u}_p(s; \mathcal{D}) - \hat{\mathbf{u}}_p(s; X)| \leq c \sqrt{\frac{\ln(\frac{2|I|}{\delta})}{2m}}.$$

Remark Given a game, we state all theorems and algorithms for an arbitrary index set I . Taking $I = P \times S$, we bound $\|\mathbf{u}(\cdot; \mathcal{D}) - \hat{\mathbf{u}}(\cdot; X)\|_\infty$.

Consider $X_{1:m}$ i.i.d. random variables, and their mean \bar{X} . (In our case, $X_j = \mathbf{u}(s; x_j)$.) Hoeffding’s inequality for bounded random variables can be used to obtain tail bounds on the probability that an empirical mean differs from its expectation. One way to characterize such bounds is to compare them to well-studied cases of common

random variables. We focus on the case of upper tails, as we apply only symmetric two-tail and upper tail bounds in this work.

Random variables that obey the Gaussian Chernoff bound can be characterized as σ_N^2 -sub-Gaussian: i.e.,

$$\mathbb{P}(\bar{X} \geq \mathbb{E}[\bar{X}] + \varepsilon) \leq \exp\left(\frac{-m\varepsilon^2}{2\sigma_N^2}\right); \text{ equivalently,}$$

$$\mathbb{P}\left(\bar{X} \geq \mathbb{E}[\bar{X}] + \sqrt{\frac{2\sigma_N^2 \ln(\frac{1}{\delta})}{m}}\right) \leq \delta.$$

Here, σ_N^2 is deemed a *variance proxy*. Using this characterization, Hoeffding's inequality reads "If X_i has range c , then X_i is $\frac{c^2}{4}$ -sub-Gaussian," because, by Popoviciu's inequality [1935], $\mathbb{V}[X_i] \leq \frac{c^2}{4}$. Thus, Hoeffding's inequality yields sub-Gaussian tail bounds.

This result is not entirely satisfying, however, as it is stated in terms of the *worst-case* variance; when $\mathbb{V}[X_i] \ll \frac{c^2}{4}$, tighter bounds should be possible. We might hope that knowledge of the variance σ^2 would imply σ^2 -sub-Gaussian, but it does not; taking the range c to ∞ allows X_i to exhibit arbitrary tail behaviors. A (σ_1^2, c) -sub-gamma [6] random variable obeys

$$\mathbb{P}\left(\bar{X} \geq \mathbb{E}[\bar{X}] + \frac{c \ln(\frac{1}{\delta})}{3m} + \sqrt{\frac{2\sigma_1^2 \ln(\frac{1}{\delta})}{m}}\right) \leq \delta.$$

While the form of such bounds is more complicated than in the sub-Gaussian case, there is an intuitive interpretation of the tail behavior as sample size increases. The key is to observe that the additive error consists of a *hyperbolic* $\frac{c \ln(\frac{1}{\delta})}{3m}$ term and a *root-hyperbolic* $\sqrt{2\sigma_1^2 \ln(\frac{1}{\delta})/m}$ term, which in learning theory are often called *fast* and *slow* terms, respectively. Sub-gamma random variables then yield *mixed convergence rates*, which initially decay quickly while the c -term dominates, before slowing to the root-hyperbolic rate when the σ_1^2 term comes to dominate.

Bennett's inequality [1962], while usually stated as a sub-Poisson bound, immediately implies that if X_i has range c and variance σ^2 , then X_i is (σ^2, c) -sub-gamma. Bennett's inequality, however, assumes the range and variance are *known*. Various *empirical* Bennett bounds have been shown [3, 4, 18], which all essentially operate by bounding the variance of a random variable in terms of its empirical variance and range, and then applying Bennett's inequality. Simply put, Bennett's inequality gives us Gaussian-like tail bounds, where the scale-dependent term acts as an *asymptotically negligible* correction for working with non-Gaussian distributions and *empirical estimates* of variance. Asymptotic central-limit-theorem bounds behave similarly, but lack corresponding finite-sample guarantees.

Our work differs from previous applications in that we require confidence intervals of *uniform width*, and thus our bounds are limited by the *maximum* variance over all parameters being estimated. The maximum variance over a set of random variables is known as the *wimpy variance* [6]. We apply a sub-gamma bound to the wimpy variance in terms of an empirical estimate, and then apply Bennett's inequality to the upper and lower tails of each utility function using the wimpy variance bound, as this, by definition, upper bounds the variance of each utility. This strategy yields uniform-width confidence intervals over all utilities, and requires a union bound over m upper tails, m lower tails, and the wimpy variance, whereas bounding each variance separately would require a union bound over m upper, m lower, and m variance bounds. Thus, our method outperforms such strategies, yielding tighter confidence intervals.

THEOREM 3.2 (BENNETT-TYPE VARIANCE-SENSITIVE FINITE-SAMPLE BOUNDS). *Suppose as in Thm. 3.1. Take*

$$\hat{v} \doteq \sup_{(p,s) \in I} \frac{1}{m-1} \sum_{j=1}^m (\mathbf{u}_p(s; x_j) - \hat{\mathbf{u}}_p(s; X))^2;$$

$$\varepsilon_v \doteq \frac{c \ln(\frac{3}{\delta})}{m-1} + \sqrt{\left(\frac{c \ln(\frac{3}{\delta})}{m-1}\right)^2 + \frac{2\hat{v} \ln(\frac{3}{\delta})}{m-1}}; \text{ \&}$$

$$\varepsilon_\mu \doteq \min \left(\underbrace{c \sqrt{\frac{\ln(\frac{3|I|}{\delta})}{2m}}}_{\text{HOEFFDING}}, \underbrace{\frac{c \ln(\frac{3|I|}{\delta})}{3m} + \sqrt{\frac{2(\hat{v} + \varepsilon_v) \ln(\frac{3|I|}{\delta})}{m}}}_{\text{EMPIRICAL BENNETT}} \right).$$

Then, with probability at least $1 - \delta$ over $X \sim \mathcal{D}^m$, it holds that

$$\sup_{(p,s) \in I} |\mathbf{u}_p(s; \mathcal{D}) - \hat{\mathbf{u}}_p(s; X)| \leq \varepsilon_\mu.$$

From the definition of ε_μ , we see that when \hat{v} is $\approx \frac{c^2}{4}$ (near-maximal), the HOEFFDING term applies, so this bound matches Theorem 3.1 to within constant factors (in particular, $\ln(\frac{3|I|}{\delta})$ instead of $\ln(\frac{2|I|}{\delta})$). On the other hand, when \hat{v} is small, Theorem 3.2 is much sharper than Theorem 3.1. A few simplifying inequalities yield

$$\varepsilon_\mu \leq \frac{7c \ln(\frac{3|I|}{\delta})}{3(m-1)} + \sqrt{\frac{2\hat{v} \ln(\frac{3|I|}{\delta})}{m}},$$

which matches the standard sub-gamma Bennett's inequality up to constant factors, with dependence on \hat{v} instead of v . In the extreme, when $\hat{v} \approx 0$ (i.e., the game is *near-deterministic*), then Theorem 3.2 improves asymptotically over Theorem 3.1 by a $\Theta(\sqrt{\ln(\frac{|I|}{\delta})/m})$ factor.

4 LEARNING ALGORITHMS

We are now ready to present our algorithms. Specifically, we discuss two Monte-Carlo sampling-based algorithms that can be used to uniformly learn empirical games, and hence ensure that the equilibria of the games they are learning are accurately approximated with high probability. Note that our algorithms apply only to finite games, as they require an enumeration of the index set I .

A conditional normal form game Γ_X , together with a black box from which we can sample distribution \mathcal{D} , serves as our mathematical model of a black-box simulator from which the utilities of a simulation-based game can be sampled. Given strategy profile s , we assume the simulator outputs a sample $\mathbf{u}_p(s, x)$, for all agents $p \in P$, after drawing a *single* condition value $x \sim \mathcal{D}$.

Our first algorithm, *global sampling* (GS), is a straightforward application of Thms. 3.1 and 3.2. The second, *progressive sampling with pruning* (PSP), iteratively prunes strategies, and thereby has the potential to expedite learning by obtaining tighter bounds than GS, given the same number of samples. We explore PSP's potential savings in our experiments (Sec. 5).

Our first algorithm, GS (Alg. 1), samples all utilities of interest, given a sample size m and a failure probability δ , and returns the ensuing empirical game together with an $\hat{\varepsilon}$ determined by either Thm. 3.1 or 3.2 that guarantees an $\hat{\varepsilon}$ -uniform approximation.

More specifically, GS takes in a conditional game Γ_X , a black box from which we can sample distribution \mathcal{D} , an index set $I \subseteq P \times S$, a sample size m , a utility range c such that utilities are required to lie in $[-c/2, c/2]$, and a bound type BD, and then draws m samples to produce an empirical game $\hat{\Gamma}_X$, represented by $\hat{\mathbf{u}}(\cdot)$, as well as an additive error $\hat{\varepsilon}$, with the following guarantee:

Algorithm 1 Global Sampling

```

1: procedure GS( $\Gamma_X, \mathcal{D}, I, m, \delta, c, \text{BD}$ )  $\rightarrow (\tilde{\mathbf{u}}, \hat{\varepsilon})$ 
2:   input: Conditional game  $\Gamma_X$ , black box from which we
   can sample distribution  $\mathcal{D}$ , index set  $I$ , sample size  $m$ , failure
   probability  $\delta$ , utility range  $c$ , and bound type  $\text{BD}$ .
3:   output: Empirical utilities  $\tilde{\mathbf{u}}, \forall (p, s) \in I$ ; additive error  $\hat{\varepsilon}$ .
4:    $\mathbf{X} \sim \mathcal{D}^m$   $\triangleright$  Draw  $m$  samples from distribution  $\mathcal{D}$ 
5:    $\forall (p, s) \in I : \tilde{\mathbf{u}}_p(s) \leftarrow \hat{\mathbf{u}}_p(s; \mathbf{X})$ 
6:   if  $\text{BD} = \text{H}$  then  $\triangleright$  See Thm. 3.1 (Hoeffding)
7:      $\hat{\varepsilon} \leftarrow c\sqrt{\frac{\ln(\frac{2|I|}{\delta})}{2m}}$ 
8:   else if  $\text{BD} = \text{B}$  then  $\triangleright$  See Thm. 3.2 (Empirical Bennett)
9:      $\hat{\varepsilon} \leftarrow \sup_{(p,s) \in I} \frac{1}{m-1} \sum_{j=1}^m (\mathbf{u}_p(s; \mathbf{x}_j) - \tilde{\mathbf{u}}_p(s))^2$ 
10:     $\varepsilon_v \leftarrow \frac{c \ln(\frac{3}{\delta})}{m-1} + \sqrt{\left(\frac{c \ln(\frac{3}{\delta})}{m-1}\right)^2 + \frac{2\hat{\varepsilon} \ln(\frac{3}{\delta})}{m-1}}$ 
11:     $\hat{\varepsilon} \leftarrow \min \left( c\sqrt{\frac{\ln(\frac{3|I|}{\delta})}{2m}}, \frac{c \ln(\frac{3|I|}{\delta})}{3m} + \sqrt{\frac{2(\hat{\varepsilon} + \varepsilon_v) \ln(\frac{3|I|}{\delta})}{m}} \right)$ 
12:   end if
13:   return  $(\tilde{\mathbf{u}}, \hat{\varepsilon})$ 
14: end procedure

```

THEOREM 4.1 (APPROXIMATION GUARANTEES OF GLOBAL SAMPLING). Consider conditional game Γ_X together with distribution \mathcal{D} and take index set $I \subseteq P \times S$ such that for all $x \in X$ and $(p, s) \in I$, $\mathbf{u}_p(s; x) \in [-c/2, c/2]$, for some $c \in \mathbb{R}$. If $\text{GS}(\Gamma_X, \mathcal{D}, I, m, \delta, c, \text{BD})$ outputs the pair $(\tilde{\mathbf{u}}, \hat{\varepsilon})$, then with probability at least $1 - \delta$, it holds that

$$\sup_{(p,s) \in I} |\mathbf{u}_p(s; \mathcal{D}) - \tilde{\mathbf{u}}_p(s)| \leq \hat{\varepsilon}.$$

Next, we present PSP (Alg. 2), which, using GS as a subroutine, draws progressively larger samples, refining the empirical game at each iteration, and stopping when the equilibria are approximated to the desired accuracy, or when the sampling budget is exhausted. Although performance ultimately depends on a game's structure, PSP can learn equilibria using vastly fewer resources than GS.

As the name suggests, PSP is a pruning algorithm. The key idea is to prune (i.e., cease estimating the utilities of) strategy profiles that (w.h.p.) are provably not equilibria. Recall that $s \in E_\varepsilon(\Gamma)$ iff $\text{Reg}_p(\Gamma, s) \leq \varepsilon$, for all $p \in P$. Thus, if there exists $p \in P$ s.t. $\text{Reg}_p(\Gamma, s) > \varepsilon$, then $s \notin E_\varepsilon(\Gamma)$. In the search for pure equilibria, such strategy profiles can be pruned.

A strategy $s \in S_p$ is said to ε -dominate another strategy $s' \in S_p$ if for all $s \in S$, taking $s' = (s_1, \dots, s_{p-1}, s'_p, s_{p+1}, \dots, s_{|P|})$, it holds that $\mathbf{u}_p(s) - \varepsilon \geq \mathbf{u}_p(s')$. The ε -rationalizable strategies $\text{Rat}_\varepsilon(\Gamma)$ are those that remain after iteratively removing all ε -dominated strategies. Only strategies in $\text{Rat}_\varepsilon(\Gamma)$ can have nonzero weight in a mixed ε -Nash equilibrium [10]; thus eliminating strategies not in $\text{Rat}_\varepsilon(\Gamma)$ is a natural pruning criterion for mixed equilibria.

If a strategy $s \in S_p$ is ε -dominated by another strategy $s' \in S_p$, then p always regrets playing strategy s , regardless of other agents' strategies. Consequently, the mixed pruning criterion is more conservative than the pure, which means more pruning occurs when learning pure equilibria.

Like GS, PSP takes in a conditional game Γ_X , a black box from which we can sample distribution \mathcal{D} , a utility range c , and a bound

type BD . Instead of a single sample size, however, it takes in a *sampling schedule* \mathbf{M} in the form of a (possibly infinite) strictly increasing sequence of integers; and instead of a single failure probability, it takes in a *failure probability schedule* δ , with each δ_t in this sequence and their sum in $(0, 1)$. These two schedules dictate the number of samples to draw and the failure probability to use at each iteration. PSP also takes in a boolean PURE that indicates whether the equilibria of interest are pure or mixed, and an *error threshold* ε , which enables early termination as soon as equilibria are estimated to within the additive factor ε .

Algorithm 2 Progressive Sampling with Pruning

```

1: procedure PSP( $\Gamma_X, \mathcal{D}, \mathbf{M}, \delta, c, \text{BD}, \text{PURE}, \varepsilon$ )  $\rightarrow ((\tilde{\mathbf{u}}, \tilde{\varepsilon}), (\hat{E}, \hat{\varepsilon}), \hat{\delta})$ 
2:   input: Conditional game  $\Gamma_X$ , black box from which we can
   sample distribution  $\mathcal{D}$ , sampling schedule  $\mathbf{M}$ , failure probability
   schedule  $\delta$ , utility range  $c$ , bound type  $\text{BD}$ , equilibrium type
    $\text{PURE}$ , error threshold  $\varepsilon$ .
3:   output: Empirical utilities  $\tilde{\mathbf{u}}, \forall (p, s) \in P \times S$ , utility error  $\tilde{\varepsilon}$ ,
   empirical equilibria  $\hat{E}$ , equilibria error  $\hat{\varepsilon}$ , failure probability  $\hat{\delta}$ .
4:    $I \leftarrow P \times S$   $\triangleright$  Initialize index set
5:    $\forall (p, s) \in I : (\tilde{\mathbf{u}}_p(s), \tilde{\varepsilon}_p(s)) \leftarrow (0, c/2)$ ,  $\triangleright$  Initialize outputs
6:   for  $t = 1, \dots, |\mathbf{M}|$  do
7:      $(\tilde{\mathbf{u}}, \tilde{\varepsilon}) \leftarrow \text{GS}(\Gamma_X, \mathcal{D}, I, \mathbf{M}_t, \delta_t, c, \text{BD})$   $\triangleright$  Improve estimates
8:      $\forall (p, s) \in I : \tilde{\varepsilon}_p(s) \leftarrow \tilde{\varepsilon}$   $\triangleright$  Update confidence intervals
9:     if  $\hat{\varepsilon} \leq \varepsilon$  or  $t = |\mathbf{M}|$  then  $\triangleright$  Termination condition
10:       $\hat{E} \leftarrow \begin{cases} \text{PURE} & : E_{2\tilde{\varepsilon}}(\tilde{\mathbf{u}}) \\ \neg \text{PURE} & : E_{2\tilde{\varepsilon}}^\circ(\tilde{\mathbf{u}}) \end{cases}$ 
11:      return  $((\tilde{\mathbf{u}}, \tilde{\varepsilon}), (\hat{E}, \hat{\varepsilon}), \sum_{i=1}^t \delta_i)$ 
12:   end if
13:    $I \leftarrow \begin{cases} \text{PURE} : \{(p, s) \in I \mid \text{Reg}_p(\tilde{\mathbf{u}}, s) \leq 2\tilde{\varepsilon}\} \\ \neg \text{PURE} : \{(p, s) \in I \mid \forall q \in P : s_q \in \text{Rat}_{2\tilde{\varepsilon}}(\tilde{\mathbf{u}})\} \end{cases}$ 
14:   end for
15: end procedure

```

In the PSP pseudocode, and in the following theorem, we overload the Reg and E operators (both pure and mixed) to depend on a utility function, rather than a game.

THEOREM 4.2. Suppose conditional game Γ_X and distribution \mathcal{D} such that for all $x \in X$ and $(p, s) \in P \times S$, $\mathbf{u}_p(s; x) \in [-c/2, c/2]$ for some $c \in \mathbb{R}$. If $\text{PSP}(\Gamma_X, \mathcal{D}, \mathbf{M}, \delta, c, \text{BD}, \text{PURE}, \varepsilon)$ outputs $((\tilde{\mathbf{u}}, \tilde{\varepsilon}), (\hat{E}, \hat{\varepsilon}), \hat{\delta})$, it holds that:

- (1) $\hat{\delta} \leq \sum_{t \in \mathbf{M}} \delta_t, \hat{\delta} \in (0, 1)$;
- (2) If $\lim_{t \rightarrow \infty} \ln(1/\delta_t)/\mathbf{M}_t = 0$, then $\hat{\varepsilon} \leq \varepsilon$.

Furthermore, if PSP terminates, then with probability at least $1 - \hat{\delta}$, the following hold simultaneously:

- (3) $|\mathbf{u}_p(s; \mathcal{D}) - \tilde{\mathbf{u}}_p(s)| \leq \tilde{\varepsilon}_p(s)$, for all $(p, s) \in P \times S$;
- (4) If PURE , then $E(\mathbf{u}) \subseteq E_{2\tilde{\varepsilon}}(\tilde{\mathbf{u}}) \subseteq E_{4\tilde{\varepsilon}}(\mathbf{u})$;
- (5) If $\neg \text{PURE}$, then $E^\circ(\mathbf{u}) \subseteq E_{2\tilde{\varepsilon}}^\circ(\tilde{\mathbf{u}}) \subseteq E_{4\tilde{\varepsilon}}^\circ(\mathbf{u})$.

Finally, we propose two possible sampling and failure probability schedules for PSP, \mathbf{M} and δ , depending on whether the sampling budget is finite or infinite. Given a finite sampling budget $m < \infty$, a neutral choice is to take \mathbf{M} to be a doubling sequence such that $\sum_{M_i \in \mathbf{M}} M_i \leq m$, with M_1 sufficiently large so as to possibly permit pruning after the first iteration (iterations that neither prune nor achieve ε -accuracy are effectively wasted), and to take $\delta_t =$

$\delta/|M|$, where δ is some maximum tolerable failure probability. This strategy may fail to produce the desired ε -approximation, as it may exhaust the sampling budget first. To guarantee a particular ε - δ -approximation, then we can take M to be an infinite doubling sequence, and δ to be a geometrically decreasing sequence such that $\sum_{t=1}^{\infty} \delta_t = \delta$, for which the conditions of Thm. 4.2 (2) hold.

At first glance, it may seem that the sample complexity of estimating a game with our empirical Bennett-type bound should depend on the number, the worst-case variance over, and the range of utility values. However, if the variance of the utilities at all equilibria is small, and if the variance of the utilities at non-equilibrium strategy profiles is no more than half the square of their regret, then PSP runs only until the high-variance high-regret strategy profiles are pruned. Once this critical point is reached, the game is essentially learned, as each subsequent iteration prunes more strategy profiles, further reducing the wimpy variance, yielding tighter bounds for the remaining profiles. Under these circumstances, PSP can provably require fewer samples than GS.

5 EXPERIMENTS

We now set out to evaluate the strength of our methodology to learn black-box games and their equilibria from samples. The empirical performance of an algorithm can vary dramatically under different distributions of inputs; in particular, the success of game-theoretic solvers can vary dramatically even within the same class of games [17, 21]. Consequently, we employ GAMUT [21], a state-of-the-art suite of game generators that is capable of producing a wide variety of interesting game inputs of varying scales with rich strategic structure, thereby affording us an opportunity to conduct a robust evaluation of our methodology. Furthermore, we employ Gambit [19], a state-of-the-art equilibrium solver. We bundled both of these packages together with our statistical learning algorithms in a python library for empirical game-theoretic analysis, pySEGTA, to make it easier for other EGTA researchers to benchmark their algorithms against ours.

Black-Box Game Design. In all our experiments, we use GAMUT to generate what we call *ground-truth games*. Ground-truth games are ordinarily inaccessible; however, we rely on them here to measure the loss experienced by our algorithms: i.e., the regrets in a learned game as compared those in the corresponding ground-truth game. To simulate a black-box game, we simply add noise drawn from a zero-centered distribution to the utilities of a ground-truth game. We detail this construction presently.

Let Γ be a realization of a ground-truth game drawn from GAMUT, and let $u_p(s)$ be the utility of player p at profile s in Γ . Fix a condition set $\mathcal{X} = [a, b]$, where $a < b$. In the conditional game $\Gamma_{\mathcal{X}}$, $u_p(s; x_{p,s}) = u_p(s) + x_{p,s}$, for $x_{p,s} \in \mathcal{X}$. Conditional game $\Gamma_{\mathcal{X}}$ together with distribution \mathcal{D} on \mathcal{X} is then our model for a black-box game. For simplicity, all noise $x_{p,s} \sim \mathcal{D}$ is drawn i.i.d..

We only consider noise distributions where \mathcal{D} is zero-centered. Consequently, the expected-normal form game $\Gamma_{\mathcal{D}}$, which is the object of our algorithms' estimation, exactly coincides with Γ : i.e., it holds that for every p and s , $u_p(s; \mathcal{D}) = \mathbb{E}_{x_{p,s} \sim \mathcal{D}} [u_p(s; x_{p,s})] = \mathbb{E}_{x_{p,s} \sim \mathcal{D}} [u_p(s) + x_{p,s}] = u_p(s) + \mathbb{E}_{x_{p,s} \sim \mathcal{D}} [x_{p,s}] = u_p(s)$, because $u_p(s)$ is constant and \mathcal{D} is zero-centered.

Experimental Setup. We normalize the utilities generated by GAMUT to lie in the range $[-10, 10]$. We experiment with three different noise regimes, high, medium, and low variance. Letting $U[a, b]$ be a uniform distribution over $[a, b]$, we model high, medium, and low variance noise by distributions $U[-2.5, 2.5]$, $U[-.5, .5]$, and $U[-.1, .1]$, respectively.

We test both GS, Alg. 1, and PSP, Alg. 2. These algorithms take as input a flag $Bd \in \{H, B\}$, indicating which bound, Hoeffding's (Thm. 3.1) or our empirical Bennett-type (Thm. 3.2), to use. Henceforth, to refer to an algorithm that uses bound Bd , we write $GS(Bd)$ and $PSP(Bd)$. Throughout our experiments, we fix $\delta = 0.05$.

Sample Efficiency of GS. In this experiment, we investigate the sampling efficiency of our algorithms; that is, the quality of the games learned, as measured by ε , as a function of the number of samples needed to achieve said guarantee. We tested ten different classes of GAMUT games, all of them two-player, with varying numbers of strategies, either indicated in parentheses next to the game's name, or two by default. For each class of games, we draw 60 random ground-truth games from GAMUT, and for each such draw, we run GS 20 times for each of sample sizes $m \in \{10, 20, 40, 80, 160, 320, 640, 1280, 2560, 5120\}$, measuring ε for all possible combinations of these parameters. We then average the measured values of ε , fixing the number of samples. Fig. 1 plots, in a log-log scale, these averages, comparing the performance of GS given Hoeffding's bound and our empirical Bennett-type bound, for the cases of high and low variance. Fig. 1 depicts only four classes of games, but trends over all ten classes are similar (see the supplemental material). In all cases, we found that our empirical Bennett-type bound produces better estimates for the same number of samples, as measured by ε . Note the initial $1/m$ decay rate, later slowing to $1/\sqrt{m}$, in the Bennett bounds, reflecting the fast c term and slow σ^2 terms of the sub-gamma bounds.

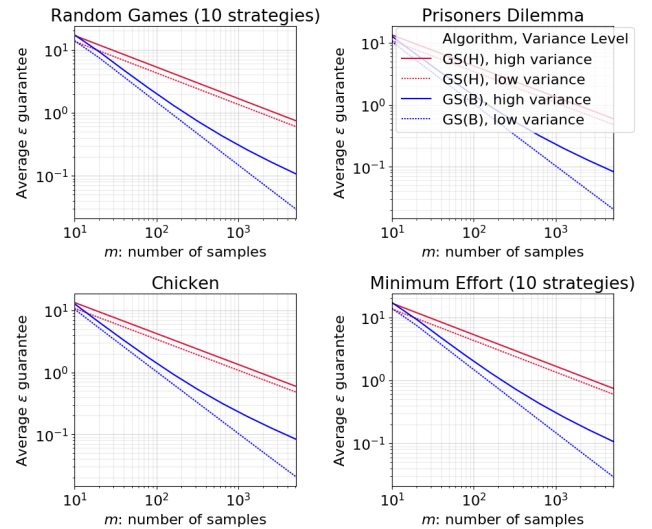


Figure 1: Quality of Learned Games

Empirical Regret of GS. In this experiment, we investigate the quality of the equilibria learned by our algorithms. To compute the equilibria of a game, we use Gambit [19], a state-of-the-art equilibrium solver. The goal of these experiments is to provide empirical evidence for Thm. 2.2, namely that our algorithms are capable of learning games that approximately preserve the Nash equilibria of black-box games. The goal is not to test the quality of different equilibrium solvers—we refer the reader to [21] for an evaluation along those lines. Hence, we fix one such solver throughout, namely, Gambit’s GNM solver, a global Newton method that computes Nash equilibria [11].

To measure the quality of learned equilibria, given a game Γ and a subset of the strategy profile space $S' \subseteq S^\circ$, we define the metric $\text{MAX-REGRET}(\Gamma, S') = \sup_{s \in S'} \max_{p \in P} \text{Reg}_p^\diamond(\Gamma, s)$: i.e., the maximum regret of any player p at any profile s in S' . Note that, given two compatible games Γ and Γ' , and S' , we can measure MAX-REGRET of *either* game, since the strategy profile space is shared by compatible games. This is useful because, given a ground-truth game Γ and a corresponding empirical estimate $\hat{\Gamma}_X$, we can measure the maximum regret of a set of Nash equilibrium profiles in Γ , say S^* , in its empirical estimate $\hat{\Gamma}_X$. Thm. 2.2 implies that, given an ε -uniform approximation $\hat{\Gamma}_X$ of Γ , we should observe $\text{MAX-REGRET}(\hat{\Gamma}_X, S^*) \leq 2\varepsilon$. Thm. 4.1 then implies that if said ε -uniform approximation holds with probability $1 - \delta$, then we should likewise observe $\text{MAX-REGRET}(\hat{\Gamma}_X, S^*) \leq 2\varepsilon$ with probability $1 - \delta$.

We empirically measure MAX-REGRET , where equilibria are computed using Gambit, for the same ten different classes of games as in the previous experiment, again over 60 draws for each class, where for each such draw, we run GS 10 times for each of sample sizes $m \in \{10, 20, 40, 80, 160\}$, measuring MAX-REGRET for all possible combinations of these parameters. We then average the measured values of MAX-REGRET , fixing the number of samples. Fig. 2 plots, in a log-log scale, both these averages (the markers) and the theoretical guarantees (the lines). This plot complements our theory, establishing experimentally that our algorithms are capable of preserving equilibria in learned black-box games. This learning is robust to various classes of games, for example, for dominance-solvable games (such as Prisoners’ Dilemma), games with guaranteed pure-strategy equilibria (such as congestion games), as well as other games with no guarantee on their equilibria other than existence (such as random and zero-sum games).³ This learning is also robust to different levels of noise, with our algorithms consistently achieving higher accuracy in practice than in theory, all across the board.

Sample Efficiency of PSP. In this experiment, we investigate the sample efficiency of PSP as compared to GS. We say that algorithm A has better *sample efficiency* than algorithm B if A requires fewer samples than B to achieve a desired accuracy ε .

Our experimental design is as follows. Fixing a game, and the following values of $\varepsilon \in \{0.125, 0.25, 0.5, 1.0\}$, we compute the number of samples $m(\varepsilon)$ that would be required for $\text{GS}(H)$ to achieve accuracy ε . We then run both $\text{GS}(H)$ and $\text{GS}(B)$ with $m(\varepsilon)$ samples.

For PSP, we use the following doubling strategy as a sampling schedule $M(m(\varepsilon)) = [m(\varepsilon)/4, m(\varepsilon)/2, m(\varepsilon), 2m(\varepsilon)]$, rounding to the nearest integer as necessary. For δ , we use a uniform schedule such

³Similar to before, we report on only three game classes here, and refer the reader to the supplemental material for results on all 10 classes.

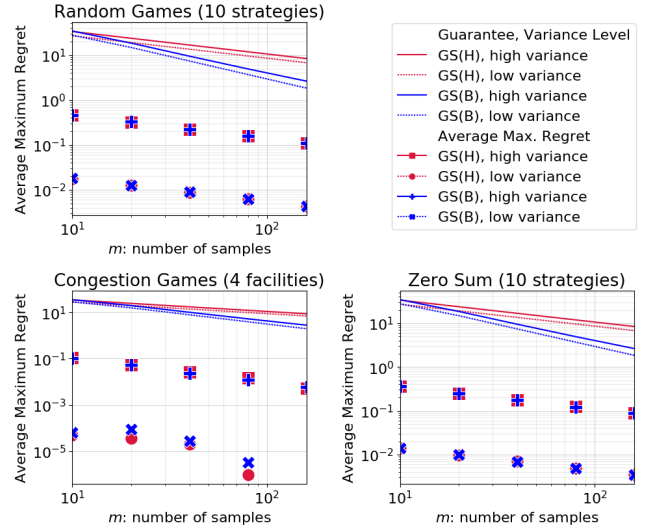


Figure 2: Average Maximum Regret

that $\sum_{\delta_t \in \delta} \delta_t = \delta$: i.e., $\delta = [0.0125, 0.0125, 0.0125, 0.0125]$. Using these schedules, we run both $\text{PSP}(H)$ and $\text{PSP}(B)$ until completion by setting the desired accuracy to zero. We prune using the mixed-strategy criterion, namely the set of rationalizable strategies.

We ran this experiment on three different classes of games: congestion games (2 players, and 2, 3, 4, and 5 facilities; game sizes 18, 98, 450, and 1,922 respectively), random games (2 players, and 5, 10, 20, and 30 strategies each; game sizes 50, 200, 800, and 1,800 respectively), and zero-sum games (2 players, and 5, 10, 20, and 30 strategies; game sizes 50, 200, 800, and 1,800 respectively). As with our other experiments, we draw multiple games for each class of games (in this case 30) and multiple runs (in this case 10 for each draw of each game). We consider medium variance only.

For all algorithms, we measure the total number of samples across all players and strategy profiles. If and when it prunes, PSP requires progressively fewer and fewer samples, with the number decreasing each iteration to the size of the unpruned game.

Table 1 summarizes the results of these experiments for select games. In all cases, we simply report the total number of samples, averaged across all experiments. The theory tells us that given this number of samples, GS must achieve at least the desired $\varepsilon \in \{0.125, 0.25, 0.5, 1.0\}$. Although there is no such guarantee for PSP , in these experiments PSP always achieved a strictly greater accuracy than GS (these accuracies are also reported in Table 1, under the columns labeled ε_{PSP}). Moreover, PSP tends to exhibit significantly better sample efficiency than GS ; notable exceptions include cases where either the games are small or the ε guarantee is loose (e.g., $\varepsilon \leq 1.0$). These results demonstrate the promise of PSP as an algorithm for learning black-box games, as its sample efficiency generally exceeds that of GS .

Limitations of PSP. While our experiments demonstrate that PSP can yield substantial savings when learning games in many different classes, in some GAMUT games, our simple doubling schedule yielded no such gains. We found Grab The Dollar to be a particularly difficult game. In Grab The Dollar, there is a prize

Bound	$\epsilon \leq 0.125$		$\epsilon \leq 0.25$		$\epsilon \leq 0.5$		$\epsilon \leq 1.0$	
	Hoeffding	Emp. Bennett	Hoeffding	Emp. Bennett	Hoeffding	Emp. Bennett	Hoeffding	Emp. Bennett
Game/Algorithm	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}	GS; PSP; ϵ_{PSP}
Congestion Games (5 facilities)	3,051; 1,654 ; 0.08	3,051; 1,449 ; 0.00	762; 464 ; 0.17	762; 364 ; 0.01	190; 146 ; 0.34	190; 93 ; 0.01	47; 58; 0.70	47; 25 ; 0.04
Zero-Sum Games (30 strategies)	2,841; 1,691 ; 0.08	2,841; 1,383 ; 0.00	710; 502 ; 0.17	710; 349 ; 0.01	177; 166 ; 0.35	177; 90 ; 0.01	44 ; 62; 0.71	44; 25 ; 0.04
Random Games (30 strategies)	2,841; 1,666 ; 0.08	2,841; 1,375 ; 0.00	710; 491 ; 0.17	710; 347 ; 0.01	177; 159 ; 0.35	177; 90 ; 0.01	44 ; 58; 0.71	44; 25 ; 0.04
Congestion Games (4 facilities)	622; 492 ; 0.09	622; 438 ; 0.00	156; 138 ; 0.17	156; 110 ; 0.01	39; 41; 0.35	39; 28 ; 0.01	10 ; 15; 0.71	10; 8 ; 0.04
Zero-Sum Games (20 strategies)	1,171; 829 ; 0.09	1,171; 708 ; 0.00	293; 240 ; 0.17	293; 179 ; 0.01	73; 77; 0.35	73; 46 ; 0.01	18 ; 28; 0.71	18; 13 ; 0.04
Random Games (20 strategies)	1,171; 809 ; 0.09	1,171; 698 ; 0.00	293; 232 ; 0.17	293; 176 ; 0.01	73; 73; 0.35	73; 45 ; 0.01	18 ; 25; 0.71	18; 12 ; 0.04
Congestion Games (3 facilities)	114 ; 145; 0.09	114 ; 135; 0.00	29; 40; 0.18	29; 34; 0.01	7; 12; 0.36	7; 9; 0.02	2; 4; 0.73	2; 2 ; 0.05
Zero-Sum Games (10 strategies)	254 ; 268; 0.09	254; 242 ; 0.00	63; 73; 0.18	63; 61 ; 0.01	16; 22; 0.36	16; 15 ; 0.02	4; 7; 0.73	4; 4 ; 0.05
Random Games (10 strategies)	254 ; 254 ; 0.09	254; 233 ; 0.00	63; 69; 0.18	63; 59 ; 0.01	16; 21; 0.36	16; 15 ; 0.02	4; 7; 0.72	4; 4 ; 0.05
Congestion Games (2 facilities)	17 ; 37; 0.09	17; 37; 0.00	4; 10; 0.19	4; 9; 0.01	1; 3; 0.38	1; 2; 0.02	1; 1; 0.76	1; 1 ; 0.05
Zero-Sum Games (5 strategies)	54 ; 94; 0.09	54; 89; 0.00	13; 25; 0.18	13; 22; 0.01	3; 7; 0.37	3; 6; 0.02	1; 2; 0.75	1; 1 ; 0.05
Random Games (5 strategies)	54; 83; 0.09	54; 90; 0.00	13; 22; 0.18	13; 20; 0.01	3; 6; 0.37	3; 5; 0.02	1; 2; 0.74	1; 1 ; 0.05

Table 1: PSP’s sample efficiency. Numbers of samples are reported in tens of thousands. The values in bold are smaller than their counterparts; as ϵ is fixed, they indicate the more sample efficient algorithms.

(or "dollar") that two players are free to grab at any time, and there are two utility values, one high and one low. If both players grab for the dollar at the same time, it will rip; so the players earn the low utility. If one grabs the dollar before the other, then that player wins the dollar (and thus high utility), while the other player earns utility somewhere between the high and the low values.

The utility structure of this game is such that the player’s utilities are the same across many different strategy profiles—in particular, whenever one player “Grabs The Dollar” before their opponent. As a result, there are few ϵ -dominated strategies, which in turn makes pruning ineffective. PSP is most effective in cases where utilities between neighboring strategy profiles (i.e., where only one player’s strategy differs) are distinct enough that pruning is possible. Arguably, this kind of structure is common in practice, where, fixing all other players’ strategies, one player’s strategy (like defect in the Prisoners’ Dilemma) can yield very different utilities than neighboring strategies (like cooperate). Finally, our PSP algorithm does not compare favorably to the baseline in games where there are few strategies, and hence few opportunities to prune.

6 SUMMARY AND CONCLUSION

In this paper, we present and evaluate a methodology for learning games that cannot be expressed analytically. On the contrary, it is assumed that a black-box simulator is available that can be queried to obtain noisy samples of utilities. In many simulation-based games of interest, queries to the simulator are exceedingly expensive, so that the time and effort required to obtain sufficiently accurate utility estimates dwarves that of any other relevant computations, including equilibrium computations. This condition holds in meta-games like Starcraft [28], for example, where agents’ choices comprise a few high-level heuristic strategies, not intractably many low-level game-theoretic strategies. Thus, our primary concern in this paper is to limit the need for sampling, while still guaranteeing that we are estimating a game well.

Our main contributions come in two flavors. First, we derive a novel bound that can adapt to any naturally occurring variance in a game’s black-box simulator, so that fewer samples are required to obtain robust guarantees on the estimated game. Whereas Tuyls et al. [28] control the FWER by applying a Šidák correction to

Hoeffding’s bound, we improve upon their approach using an empirical variant of Bennett’s inequality that does not require *a priori* knowledge of the variance.

Second, we develop an algorithm that progressively samples a game, all the while pruning strategy profiles: i.e., ceasing to estimate those strategy profiles that provably (with high probability) do not comprise any (approximate) equilibria. In extensive experimentation over a broad swath of games, we show that this algorithm, equipped with our variance-sensitive bound, makes frugal use of samples, often requiring far fewer to learn to the same—or even a better—degree of accuracy than a variant of the sampling algorithm in Tuyls et al. [28]’s, which serves as a baseline.

A pySEGTA: A PYTHON LIBRARY

We carried out our experiments in a Python library we developed and named pySEGTA, for *statistical* EGTA.⁴ pySEGTA interfaces with both GAMUT and Gambit, exposing simple interfaces by which users can generate games (GAMUT), learn them (via our learning algorithms, for example), and solve them (Gambit). As the logic concerning game implementation is entirely separate from game learning and/or solving, pySEGTA can be used to analyze arbitrarily complex simulation-based games with arbitrarily complex strategies. pySEGTA already affords access to most GAMUT games, and is designed to be easily extensible to interface with other game generators. To do so only requires describing a game’s structure (number of players, and per-player numbers of strategies), and implementing one query method, which takes as input a strategy profile and returns sample utilities for all players at the given strategy profile. pySEGTA also includes parameterizable implementations of both GS and PSP, and was designed with extensibility in mind, so that other users can incorporate their learning algorithms as they are developed. Our intent is that pySEGTA ease the work of benchmarking empirical game-theoretic learning algorithms.

B ACKNOWLEDGEMENTS

This work was supported in part by NSF Award CMMI-1761546, NSF grant RI-1813444, and DARPA/AFRL grant FA8750-19-2-1006.

⁴pySEGTA is publicly accessible at <http://github.com/eareyan/pysegta>.

REFERENCES

- [1] Enrique Areyan Viqueira, Cyrus Cousins, Yasser Mohammad, and Amy Greenwald. 2019. Empirical Mechanism Design: Designing Mechanisms from Data. In *UAI*. AUA Press, 406.
- [2] Enrique Areyan Viqueira, Amy Greenwald, Cyrus Cousins, and Eli Upfal. 2019. Learning Simulation-Based Games from Data. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 1778–1780.
- [3] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. 2007. Tuning bandit algorithms in stochastic environments. In *International conference on algorithmic learning theory*. Springer, 150–165.
- [4] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. 2007. Variance estimates and exploration function in multi-armed bandit. In *CERTIS Research Report 07–31*. Citeseer.
- [5] George Bennett. 1962. Probability inequalities for the sum of independent random variables. *J. Amer. Statist. Assoc.* 57, 297 (1962), 33–45.
- [6] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. 2013. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- [7] Constantinos Daskalakis, Paul W Goldberg, and Christos H Papadimitriou. 2009. The complexity of computing a Nash equilibrium. *SIAM J. Comput.* 39, 1 (2009), 195–259.
- [8] Bradley Efron and Robert J. Tibshirani. 1993. *An Introduction to the Bootstrap*. Number 57 in Monographs on Statistics and Applied Probability. Chapman & Hall/CRC, Boca Raton, Florida, USA.
- [9] Tapio Elomaa and Matti Kääriäinen. 2002. Progressive Rademacher sampling. In *AAAI/IAAI*. 140–145.
- [10] Robert Gibbons. 1992. *Game theory for applied economists*. Princeton University Press.
- [11] Srihari Govindan and Robert Wilson. 2003. A global Newton method to compute Nash equilibria. *Journal of Economic Theory* 110, 1 (2003), 65–86.
- [12] Wassily Hoeffding. 1963. Probability inequalities for sums of bounded random variables. *Journal of the American statistical association* 58, 301 (1963), 13–30.
- [13] Patrick R Jordan, Christopher Kiekintveld, and Michael P Wellman. 2007. Empirical game-theoretic analysis of the TAC supply chain game. In *Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems*. ACM, 193.
- [14] Patrick R Jordan, Yevgeniy Vorobeychik, and Michael P Wellman. 2008. Searching for approximate equilibria in empirical games. In *Proceedings of the 7th international joint conference on Autonomous agents and multiagent systems-Volume 2*. International Foundation for Autonomous Agents and Multiagent Systems, 1063–1070.
- [15] Patrick R Jordan and Michael P Wellman. 2010. Designing an ad auctions game for the trading agent competition. In *Agent-Mediated Electronic Commerce. Designing Trading Strategies and Mechanisms for Electronic Markets*. Springer.
- [16] Wolfgang Ketter, Markus Peters, and John Collins. 2013. Autonomous Agents in Future Energy Markets: The 2012 Power Trading Agent Competition.. In *AAAI*.
- [17] Kevin Leyton-Brown, Eugene Nudelman, and Yoav Shoham. 2002. Learning the empirical hardness of optimization problems: The case of combinatorial auctions. In *International Conference on Principles and Practice of Constraint Programming*. Springer, 556–572.
- [18] Andreas Maurer and Massimiliano Pontil. 2009. Empirical Bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740* (2009).
- [19] McLennan Andrew M. McKelvey, Richard D. and Theodore L. Turocy. 2019. Gambit: Software Tools for Game Theory, Version 16.0.1. (2019). Retrieved November 7, 2019 from <http://www.gambit-project.org>.
- [20] John F Nash. 1950. Equilibrium points in n-person games. *Proceedings of the national academy of sciences* (1950).
- [21] Eugene Nudelman, Jennifer Wortman, Yoav Shoham, and Kevin Leyton-Brown. 2004. Run the GAMUT: A comprehensive approach to evaluating game-theoretic algorithms. In *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems-Volume 2*. IEEE Computer Society, 880–887.
- [22] Victor Picheny, Mickael Binois, and Abderrahmane Habbal. 2016. A Bayesian optimization approach to find Nash equilibria. *arXiv preprint arXiv:1611.02440* (2016).
- [23] Tiberiu Popoviciu. 1935. Sur les équations algébriques ayant toutes leurs racines réelles. *Mathematica* 9 (1935), 129–145.
- [24] Matteo Riondato and Eli Upfal. 2015. Mining frequent itemsets through progressive sampling with Rademacher averages. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1005–1014.
- [25] Matteo Riondato and Eli Upfal. 2018. ABRA: Approximating betweenness centrality in static and dynamic graphs with Rademacher averages. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 12, 5 (2018), 61.
- [26] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature* 529, 7587 (2016), 484.
- [27] Anderson Tavares, Hector Azpurua, Amanda Santos, and Luiz Chaimowicz. 2016. Rock, paper, starcraft: Strategy selection in real-time strategy games. In *The Twelfth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment (AIIDE-16)*.
- [28] Karl Tuyls, Julien Perolat, Marc Lanctot, Joel Z Leibo, and Thore Graepel. 2018. A Generalised Method for Empirical Game Theoretic Analysis. *arXiv preprint arXiv:1803.06376* (2018).
- [29] Leslie G Valiant. 1984. A theory of the learnable. *Commun. ACM* 27, 11 (1984), 1134–1142.
- [30] Enrique Areyan Viqueira, Cyrus Cousins, Eli Upfal, and Amy Greenwald. 2019. Learning Equilibria of Simulation-Based Games. *arXiv preprint arXiv:1905.13379* (2019).
- [31] Yevgeniy Vorobeychik. 2010. Probabilistic analysis of simulation-based games. *ACM Transactions on Modeling and Computer Simulation (TOMACS)* 20, 3 (2010), 16.
- [32] Yevgeniy Vorobeychik, Christopher Kiekintveld, and Michael P Wellman. 2006. Empirical mechanism design: methods, with application to a supply-chain scenario. In *Proceedings of the 7th ACM conference on Electronic commerce*. ACM.
- [33] Yevgeniy Vorobeychik and Michael P Wellman. 2008. Stochastic search methods for Nash equilibrium approximation in simulation-based games. In *Proceedings of the 7th international joint conference on Autonomous agents and multiagent systems-Volume 2*. International Foundation for Autonomous Agents and Multiagent Systems, 1055–1062.
- [34] Michael P Wellman. 2006. Methods for empirical game-theoretic analysis. In *AAAI*. 1552–1556.
- [35] Michael P Wellman, Tae Hyung Kim, and Quang Duong. 2013. Analyzing incentives for protocol compliance in complex domains: A case study of introduction-based routing. *arXiv preprint arXiv:1306.0388* (2013).
- [36] Bryce Wiedenbeck. 2014. Approximate Game Theoretic Analysis for Large Simulation-based Games. In *Proceedings of the 2014 International Conference on Autonomous Agents and Multi-agent Systems (AAMAS '14)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1745–1746. <http://dl.acm.org/citation.cfm?id=2615731.2616156>

SUPPLEMENTAL MATERIAL OF IMPROVED ALGORITHMS FOR LEARNING EQUILIBRIA OF SIMULATION-BASED GAMES, PAPER #1318

Proofs

PROOF OF LEMMA 2.1. For any agent p and mixed strategy profile $\tau \in S^\circ$, $\mathbf{u}_p(\tau) = \sum_{s \in S} \tau(s) \mathbf{u}_p(s)$, where $\tau(s) = \prod_{p' \in P} \tau_{p'}(s_{p'})$. So, $\mathbf{u}_p(\tau) - \mathbf{u}'_p(\tau) = \sum_{s \in S} \tau(s)(\mathbf{u}_p(s) - \mathbf{u}'_p(s)) \leq \sup_{s \in S} |\mathbf{u}_p(s) - \mathbf{u}'_p(s)|$, by Hölder's inequality. Hence, $\sup_{\tau \in S^\circ} |\mathbf{u}_p(\tau) - \mathbf{u}'_p(\tau)| \leq \sup_{s \in S} |\mathbf{u}_p(s) - \mathbf{u}'_p(s)|$, from which it follows that $\sup_{p \in P, \tau \in S^\circ} |\mathbf{u}_p(\tau) - \mathbf{u}'_p(\tau)| \leq \|\Gamma - \Gamma'\|_\infty$. Equality holds for any p and s that realize the supremum in $\|\Gamma - \Gamma'\|_\infty$, as any such pure strategy profile is also mixed. \square

PROOF OF THEOREM 2.2. First note the following: if $A \subseteq B$, then $C \cap A \subseteq C \cap B$. Hence, since any pure Nash equilibrium is also a mixed Nash equilibrium, taking C to be the set of all pure strategy profiles, we need only show 2. We do so by showing $E_\gamma^\circ(\Gamma) \subseteq E_{2\varepsilon+\gamma}^\circ(\Gamma')$, for $\gamma \geq 0$, which implies both containments, taking $\gamma = 0$ for the lesser, and $\gamma = 2\varepsilon$ for the greater.

Suppose $\mathbf{s} \in E_\gamma^\circ(\Gamma)$, for all $p \in P$. We will show that \mathbf{s} is $(2\varepsilon + \gamma)$ -optimal in Γ' , for all $p \in P$. Fix an agent p , and define $T_{p,s}^\circ \doteq \{\tau \in S^\circ \mid \tau_q = s_q, \forall q \neq p\}$. In words, $T_{p,s}^\circ$ is the set of all mixed strategy profiles in which the strategies of all agents $q \neq p$ are fixed at s_q . Now take $\mathbf{s}^* \in \arg \max_{\tau \in T_{p,s}^\circ} \mathbf{u}_p(\tau)$ and $\mathbf{s}'^* \in \arg \max_{\tau \in T_{p,s}^\circ} \mathbf{u}'_p(\tau)$. Then:

$$\begin{aligned} \text{Reg}_p(\Gamma', \mathbf{s}) &= \mathbf{u}'_p(\mathbf{s}'^*) - \mathbf{u}'_p(\mathbf{s}) \\ &\leq (\mathbf{u}_p(\mathbf{s}'^*) + \varepsilon) - (\mathbf{u}_p(\mathbf{s}) - \varepsilon) \\ &\leq (\mathbf{u}_p(\mathbf{s}^*) + \varepsilon) - (\mathbf{u}_p(\mathbf{s}) - \varepsilon) \\ &\leq (\mathbf{u}_p(\mathbf{s}^*) + \varepsilon) - (\mathbf{u}_p(\mathbf{s}^*) - \varepsilon - \gamma) = 2\varepsilon + \gamma \end{aligned}$$

The first line follows by definition. The second holds by Lemma 2.1 and the fact that Γ' is a uniform ε -approximation of Γ , the third as \mathbf{s}^* is optimal for p in Γ , and the fourth as \mathbf{s} is a γ -Nash in Γ . \square

PROOF OF THEOREM 4.1. Use Theorem 3.1 when BOUND = Hoeffding and Theorem 3.2 when BOUND = Bennett. \square

Variance Bound Proofs. We first show that the empirical wimpy variance can be used to upper-bound the true wimpy variance. In this section, \mathbf{x}

$$v \doteq \sup_{f \in \mathcal{F}} \mathbb{V}[f] \text{ \& \; } \hat{v} \doteq \sup_{f \in \mathcal{F}} \frac{1}{m-1} \sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}[f])^2.$$

LEMMA .1 (EMPIRICAL WIMPY VARIANCE EXPECTATION). Suppose $\mathbf{x} \sim \mathcal{D}^m$. Then $v \leq \mathbb{E}_{\mathbf{x}}[\hat{v}]$.

PROOF. We present a single inequality chain, boxing each term in the lemma statement for emphasis.

$$\begin{aligned} v &= \sup_{f \in \mathcal{F}} \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - \mathbb{E}[f])^2 \right] && \text{DEFINITION OF } v \\ &= \sup_{f \in \mathcal{F}} \mathbb{E} \left[\frac{1}{m-1} \sum_{i=1}^m (f(\mathbf{x}_i) - \hat{\mathbb{E}}_{\mathbf{x}}[f])^2 \right] && \text{BESSEL'S CORRECTION} \\ &\leq \mathbb{E}_{\mathbf{x}} \left[\sup_{f \in \mathcal{F}} \frac{1}{m-1} \sum_{i=1}^m (f(\mathbf{x}_i) - \hat{\mathbb{E}}_{\mathbf{x}}[f])^2 \right] && \text{JENSEN'S INEQUALITY} \\ &= \mathbb{E}_{\mathbf{x}}[\hat{v}] && \text{DEFINITION OF } \hat{v} \end{aligned}$$

\square

The next result employs *self-bounding functions* [6] to show that with high probability the true wimpy variance is not much larger than the empirical wimpy variance. A similar result can be found in [18] for the ordinary variance and sample variance.

LEMMA .2. Suppose $\mathcal{F} \subseteq \mathcal{X} \rightarrow [-c/2, c/2]$ and $X_{1:m} \sim \mathcal{D}^m$. Then letting

$$\varepsilon_{v-} \doteq \frac{c \ln(\frac{1}{\delta})}{m-1} + \sqrt{\left(\frac{c \ln(\frac{1}{\delta})}{m-1} \right)^2 + \frac{2 \ln(\frac{1}{\delta}) \hat{v}}{m-1}},$$

it holds that

$$\mathbb{P}(v \geq \hat{v} + \varepsilon_{v-}) \leq \delta,$$

PROOF. Suppose WLOG that $c = 1$ (the result follows via rescaling for $c \neq 1$). Suppose $X \in \mathcal{X}^m$, let $X_{\setminus j}$ denote the vector X , missing the j th component, and take $Z(X) \doteq \sup_{f \in \mathcal{F}} m \hat{\mathbb{V}}_X[f] = \sup_{f \in \mathcal{F}} \sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2$, and $Z_j(X) \doteq \sup_{f \in \mathcal{F}} (m-1) \hat{\mathbb{V}}_{X_{\setminus j}}[f] = \sup_{f \in \mathcal{F}} \sum_{i=1, i \neq j}^m (f(X_i) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2$, where $\hat{\mathbb{V}}_X[f]$ denotes the (biased) sample variance of f over X .

We first show that

$$Z_j(X) = \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} (f(X_j) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 ; \quad (1)$$

to see this, consider that

$$\begin{aligned} Z_j(X) &= \sup_{f \in \mathcal{F}} \sum_{i=1, i \neq j}^m (f(X_i) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 = \sup_{f \in \mathcal{F}} \sum_{i=1, i \neq j}^m \left((f(X_i) - \hat{\mathbb{E}}_X[f]) + (\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f]) \right)^2 && \text{ADDITIVE IDENTITY} \\ &&& \text{DEFINITION OF } Z_j \\ &= \sup_{f \in \mathcal{F}} \sum_{i=1, i \neq j}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 + 2(f(X_i) - \hat{\mathbb{E}}_X[f])(\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f]) + (\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 && \text{UNFACTORING} \\ &= \sup_{f \in \mathcal{F}} \sum_{i=1, i \neq j}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 + (2f(X_i) - \hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f])(\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f]) && \text{ALGEBRA} \\ &= \sup_{f \in \mathcal{F}} \left(\sum_{i=1, i \neq j}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) + (m-1)(\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f])(\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f]) && \text{COMMUTATIVITY} \\ &&& \text{AVERAGES} \\ &= \sup_{f \in \mathcal{F}} \left(\sum_{i=1, i \neq j}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - (m-1)(\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 && \text{ALGEBRA} \\ &= \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - (m-1)(\hat{\mathbb{E}}_X[f] - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 - (f(X_j) - \hat{\mathbb{E}}_X[f])^2 && \text{ADDITIVE IDENTITY} \\ &= \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - (m-1)\left(\frac{1}{m}f(X_j) - \frac{1}{m}\hat{\mathbb{E}}_{X_{\setminus j}}[f]\right)^2 - \left(\frac{m-1}{m}f(X_j) - \frac{m-1}{m}\hat{\mathbb{E}}_{X_{\setminus j}}[f]\right)^2 && \text{AVERAGES} \\ &= \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} (f(X_j) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 . && \forall u \in \mathbb{R} : (m-1)\left(\frac{1}{m}u\right)^2 - \left(\frac{m-1}{m}u\right)^2 = \frac{m-1}{m}u^2 \end{aligned}$$

Here AVERAGES justifications follow via arduous but ultimately uninteresting commutativity manipulations.

The desideratum reduces to showing that $\frac{m}{m-1}Z$ is an $(\frac{m}{m-1}, 0)$ -self-bounding function. *Nonnegativity* is clear from the definition of $Z(\cdot)$. We now show *underestimation*; first consider that for any $j \in \{1, \dots, m\}$,

$$\begin{aligned} Z_j(X) &= \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} (f(X_j) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 && \text{EQ. (1)} \\ &\geq \left(\sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) \right) - \left(\sup_{f \in \mathcal{F}} \frac{m-1}{m} (f(X_j) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 \right) && \text{PROPERTIES OF SUPREMA} \\ &\geq Z(X) - \frac{m-1}{m} . && \text{DEFINITION OF } Z(\cdot) \\ &&& \forall f \in \mathcal{F} : \text{image}(f) \subseteq [0, 1] \end{aligned}$$

We may thus conclude that $\frac{m}{m-1}Z_j(X) \geq \frac{m}{m-1}Z(X) - 1$, which satisfies *underestimation*. We now show that $Z(\cdot)$, and subsequently $\frac{m}{m-1}Z(\cdot)$, obey $(\frac{m}{m-1}, 0)$ -self-boundedness.

$$\begin{aligned} \sum_{j=1}^m Z(X) - Z_j(X) &= mZ(X) - \sum_{j=1}^m \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} (f(X_j) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 && \text{EQ. 1} \\ &\leq mZ(X) - \sup_{f \in \mathcal{F}} \sum_{j=1}^m \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} (f(X_j) - \hat{\mathbb{E}}_{X_{\setminus j}}[f])^2 && \text{SUBADDITIVITY} \end{aligned}$$

$$\begin{aligned}
&= mZ(X) - \sup_{f \in \mathcal{F}} m \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} \sum_{j=1}^m (f(X_j) - \hat{\mathbb{E}}_{\mathbf{x}_j}[f])^2 && \text{LINEARITY} \\
&= mZ(X) - \sup_{f \in \mathcal{F}} m \left(\sum_{i=1}^m (f(X_i) - \hat{\mathbb{E}}_X[f])^2 \right) - \frac{m-1}{m} \sum_{j=1}^m \left(\frac{m}{m-1} f(X_j) - \frac{m}{m-1} \hat{\mathbb{E}}_X[f] \right)^2 && \hat{\mathbb{E}}_{\mathbf{x}_j}[f] = \frac{1}{m-1} (m \hat{\mathbb{E}}_X[f] - f(X_j)) \\
&= \frac{m(m-1)}{m-1} Z(X) - \left(\frac{m(m-1)}{m-1} - \frac{m}{m-1} \right) Z(X) = \frac{m}{m-1} Z(X) && \text{DEFINITION OF } Z(\cdot)
\end{aligned}$$

It is now a straightforward matter to observe that $\frac{m}{m-1}Z(\cdot)$, also obeys $(\frac{m}{m-1}, 0)$ -self-boundedness, as

$$\sum_{j=1}^m Z(X) - Z_j(X) \leq \frac{m}{m-1} Z(X) \implies \frac{m}{m-1} Z(X) - \frac{m}{m-1} Z_j(X) \leq \frac{m}{m-1} \cdot \left(\frac{m}{m-1} Z(X) \right).$$

With all three desiderata shown, we may now conclude that $\frac{m}{m-1}Z(\cdot)$ is a $(\frac{m}{m-1}, 0)$ -self-bounding function. We now translate this into upper-tail bounds via self-bounding function inequalities [6]:

$$\begin{aligned}
\mathbb{P}_{\mathbf{x}}(v \geq \frac{m}{m-1} \hat{v} + \varepsilon) &\leq \mathbb{P}_{\mathbf{x}} \left(\frac{m}{m-1} \mathbb{E}[Z] \geq \frac{m}{m-1} Z(\mathbf{x}) + m\varepsilon \right) && \begin{array}{l} m\hat{v} = Z(\mathbf{x}) \\ \text{LEMMA .1} \implies v \leq \frac{1}{m-1} \mathbb{E}[Z] \end{array} \\
&\leq \exp \left(\frac{-m^2 \varepsilon^2}{2 \frac{m^2}{(m-1)^2} \mathbb{E}[Z]} \right) && \begin{array}{l} \text{LINEARITY} \\ \text{WSBF INEQUALITY} \end{array} \\
&= \exp \left(\frac{-m^2 \varepsilon^2}{\frac{2m^2}{m-1} \mathbb{E}[\hat{v}]} \right) = \exp \left(\frac{-(m-1)\varepsilon^2}{2 \mathbb{E}[\hat{v}]} \right) && \text{ALGEBRA}
\end{aligned}$$

Finally, we complete the proof by noting that, in the case where the bound holds, it also holds that $\mathbb{E}[\hat{v}] \leq \hat{v} + \varepsilon$, thus

$$\mathbb{P}_{\mathbf{x}}(v \geq \hat{v} + \varepsilon) \leq \delta, \quad \forall \varepsilon : \delta = \exp \left(\frac{-m^2 \varepsilon^2}{\frac{2m^2}{(m-1)^2} Z(\mathbf{x}) + \frac{2m^2}{m-1} \varepsilon} \right) = \exp \left(\frac{-(m-1)\varepsilon^2}{2(\hat{v} + \varepsilon)} \right),$$

which implies the result. \square

We now show a refinement of theorem 3.2, from which theorem 3.2 follows directly.

THEOREM .3. *Suppose as in theorem 3.2. Take*

$$\begin{aligned}
\varepsilon_{v-}(\delta) &\doteq \frac{c \ln(\frac{3}{\delta})}{m} + \sqrt{\left(\frac{c \ln(\frac{3}{\delta})}{m} \right)^2 + \frac{2\hat{v} \ln(\frac{3}{\delta})}{m}}, \\
\varepsilon_{\mathbb{B}} &\doteq \min \left(\underbrace{c \sqrt{\frac{\ln(\frac{3|I|}{\delta})}{2m}}}_{\text{HOEFFDING}}, \underbrace{\frac{c \ln(\frac{3|I|}{\delta})}{3m} + \sqrt{\left(\frac{c \ln(\frac{3|I|}{\delta})}{3m} \right)^2 + \frac{2 \ln(\frac{3|I|}{\delta})(\hat{v} + \varepsilon_{v-}(\frac{\delta'}{3}))}{m}}}_{\text{BENNETT}} \right) \\
&\leq \frac{4c \ln(\frac{3|I|}{\delta})}{3(m-1)} + \sqrt{\frac{2 \ln(\frac{3|I|}{\delta}) \hat{v}}{m-1}}.
\end{aligned}$$

Now, for any $\delta \in (0, 1)$, we may bound

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \mathbb{E}[f] - \hat{\mathbb{E}}_X[f] \right| \geq \varepsilon_{\mathbb{B}} \right) \leq \delta.$$

PROOF. This result holds via a union bound on the wimpy variance over all utility values, via lemma .2 with probability $\frac{\delta}{3}$, and then a union bound for Bennet's inequality on each of the $|I|$ utilities being estimated. For completeness, we reproduce the sub-Gamma form of Bennet's inequality here:

$$\mathbb{P}(|\bar{\mathbf{x}} - \mu| \geq \varepsilon) \leq 2 \exp \left(\frac{-m\varepsilon^2}{2(\sigma^2 + \frac{\varepsilon}{3}\varepsilon)} \right) \Leftrightarrow \mathbb{P} \left(|\bar{\mathbf{x}} - \mu| \geq \frac{c \ln(\frac{3}{\delta})}{3m} + \sqrt{\left(\frac{c \ln(\frac{3}{\delta})}{3m} \right)^2 + \frac{2\sigma^2 \ln(\frac{3}{\delta})}{m}} \right) \leq \delta.$$

Note that the variance of each utility is upper-bounded by the wimpy variance (by definition), which is itself bounded via lemma .2. \square

PROOF OF THEOREM 4.2. To see 1, note that $\hat{\delta}$ is computed on line 11 as a partial sum of δ_t , each addend and the sum of which are all by assumption on $(0, 1)$; thus the result holds.

To see 2, note that if $\lim_{t \rightarrow \infty} \ln(1/\delta_t)/M_t = 0$, then both the Hoeffding's and Bennett's bounds employed by GS tend to 0, as both decay asymptotically (in expectation) as $O(\sqrt{\ln(1/\delta_t)/M_t})$ (see Theorems 3.1 and 3.2). For infinite sampling schedules, the termination condition of line 9 ($\hat{\epsilon} \leq \epsilon$) is eventually met, as $\hat{\epsilon}$ is the output of GS, and thus 2 holds.

To establish 3, we show the following: assuming termination occurs at timestep n , with probability at least $1 - \hat{\delta}$, at every t in $\{1, \dots, n\}$, it holds that $\sup_{(p,s) \in P \times S} |u_p(s; \mathcal{D}) - \tilde{u}_p(s)| \leq \tilde{\epsilon}_p(s)$. This property follows from the GS guarantees of Theorem 4.1, as at each timestep t , the guarantee holds with probability at least $1 - \delta_t$; thus by a union bound, the guarantees hold simultaneously at all time steps with probability at least $1 - \sum_{i=1}^n \delta_i = 1 - \hat{\delta}$. That the GS guarantees hold for unpruned indices should be clear; for pruned indices, since only error bounds for indices updated on line 7 are tightened on line 8, it holds by the GS guarantees of previous iterations.

Without pruning, 4 and 5 would follow directly from 3 via Theorem 2.2, but with pruning, the situation is a bit more involved. To see 4, observe that at each time step, only indices (p, s) such that $\text{Reg}_p(\tilde{u}(\cdot), s) > 2\hat{\epsilon}$ are pruned (line 13), thus we may guarantee that with probability at least $1 - \hat{\delta}$, $\text{Reg}_p(\mathbf{u}(\cdot; \mathcal{D}), s) > 0$. Increasing the accuracy of the estimates of these strategy profiles is thus not necessary, as they do not comprise pure equilibria (w.h.p.), and they will never be required to refute equilibria, as these will never be a best response for any agent from any strategy profile.

5 follows similarly, except that nonzero regret implies that a pure strategy profile is not a pure Nash equilibrium, but it does *not* imply that it is *not* part of any mixed Nash equilibrium. Consequently, we use the more conservative pruning criterion of strategic dominance (a strategy is dominated if it is not rationalizable), requiring 2ϵ -dominance in $\tilde{\mathbf{u}}$, as this implies nonzero dominance in \mathbf{u} . \square

More Experimental Results

Figure 3 shows the complete set of results on the *quality of learned games* for all 10 classes of GAMUT games tested in this paper.

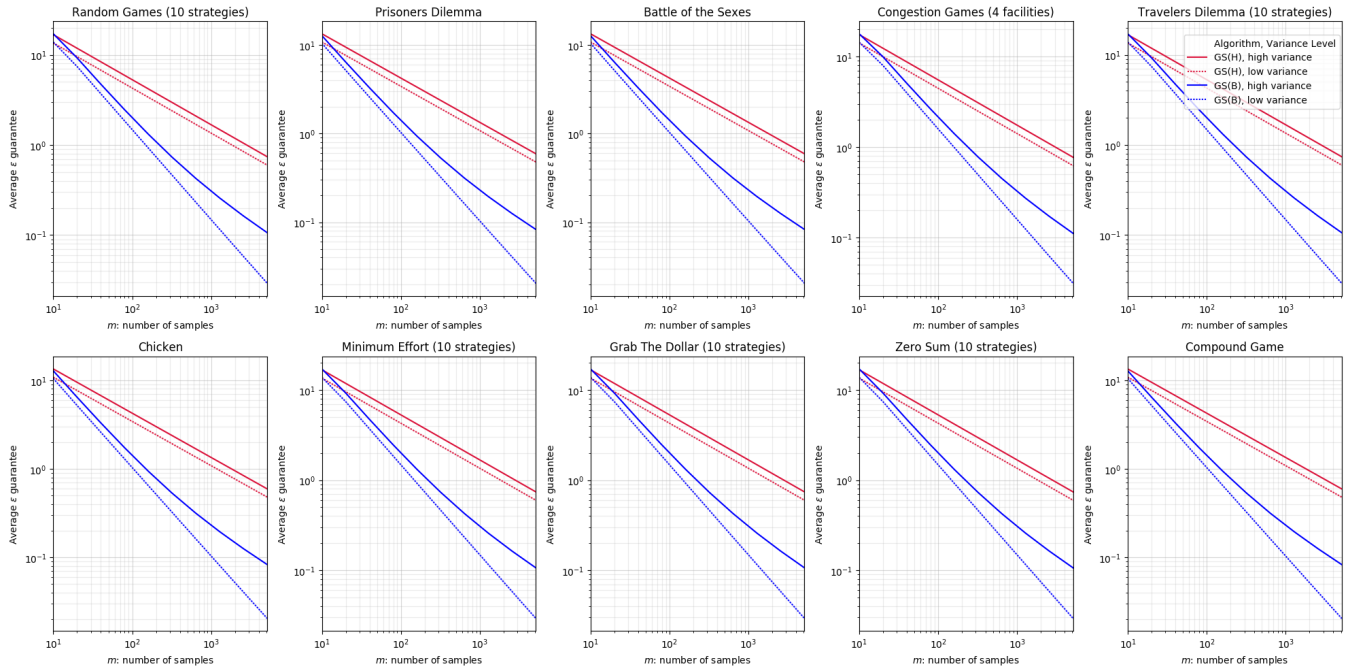


Figure 3: Quality of Learned Games - All 10 classes of games tested

Figure 4 shows the complete set of results on the *empirical regret* for all 10 classes of GAMUT games tested in this paper.

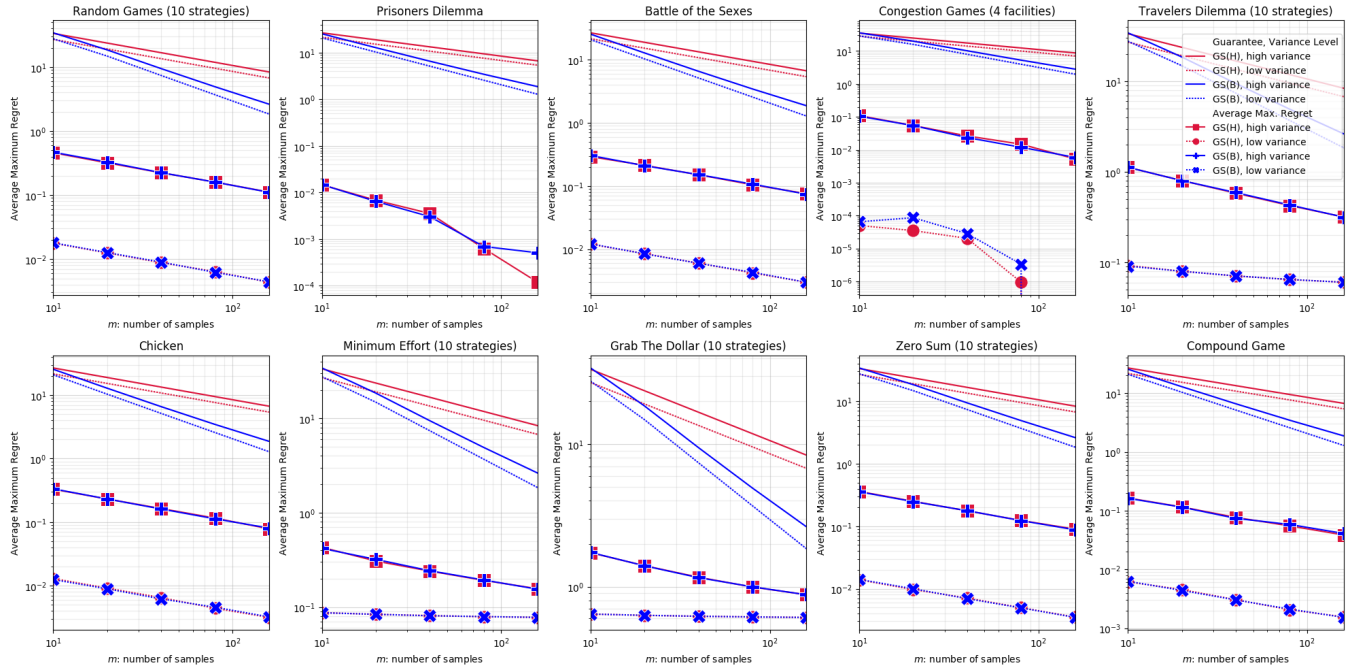


Figure 4: Average Maximum Regret - All 10 classes of games tested