

1

Solving infinite games with deep multiagent reinforcement learning

2

3

Thesis proposal

4

Carlos G. Martin

5

Computer Science Department

6

School of Computer Science

7

Carnegie Mellon University

8

Pittsburgh, PA 15213

9

Thesis committee

10

Tuomas Sandholm (CMU, chair)

11

Vincent Conitzer (CMU)

12

Fei Fang (CMU)

13

Ian Gemp (Google)

14

*Submitted in partial fulfillment of the requirements
for the degree of Doctor of Philosophy*

15

16

April 2026

Abstract

18 Game theory provides a rigorous mathematical framework for analyzing strategic interactions between
19 agents with interdependent goals. It can be used to prescribe rational behavior for individuals,
20 predict emergent dynamics in multiagent systems, and design effective social mechanisms.

21 In this thesis, we study the problem of solving infinite games. Such games can have infinitely
22 many states, actions, players, and steps. Unlike mean-field games, our players need not be symmetric
23 or exchangeable. Furthermore, we allow such games to have partial observability, hidden information,
24 imperfect recall, stochastic state transitions, discontinuous utility functions, and interdependent
25 social preferences (i.e., matrix-valued discount factors). Together, these properties can model a wide
26 range of highly complex, real-world scenarios that defy traditional game-theoretic solvers.

27 To tackle this problem, we propose a unified framework grounded in deep multiagent reinforce-
28 ment learning. It includes five core components. First, it introduces randomized policy networks
29 (RPNs) to model observation-dependent mixed strategies over infinite action spaces. Second, it
30 represents complex strategy profiles across an infinite continuum of players using player-to-strategy
31 networks (P2SNs). Third, it evolves these representations through a shared-parameter simultaneous
32 gradient (SPSG), which extends the standard simultaneous gradient to this shared-parameter regime.
33 Fourth, to ensure computational efficiency, it estimates this gradient using randomized parameter
34 perturbations via a joint-perturbation simultaneous pseudo-gradient (JPSPG). Fifth, it employs
35 approximate exploitability descent (ApproxED) with learned best-response functions (BRFs).

36 We propose to benchmark our approach on a diverse suite of real-world domains. These include
37 financial markets, traffic flow, epidemiological contagion, energy grids, and population ecology.

38 Acknowledgments

39 This material is based on work supported by the following sources of funding:

- 40 • **National Science Foundation (NSF)** grants
 - 41 ○ IIS-1901403
 - 42 ○ CCF-1733556
 - 43 ○ RI-2312342
 - 44 ○ RI-1901403
- 45 • **Army Research Office (ARO)** awards
 - 46 ○ W911NF2210266
- 47 • **Vannevar Bush Faculty Fellowship (VBFF) Office of Naval Research (ONR)** awards
 - 48 ○ N00014-23-1-2876
- 49 • **National Institutes of Health (NIH)** awards
 - 50 ○ A240108S001

51 Contents

52	1 Introduction	1
53	1.1 Background	1
54	1.2 Goal	2
55	1.3 Motivation	3
56	1.3.1 Infinite-action games	3
57	1.3.2 Infinite-player games	4
58	1.3.3 Matrix-valued discount factors	4
59	1.4 Outline	7
60	2 Mathematical preliminaries	8
61	2.1 Notation	8
62	2.2 Set theory	9
63	2.3 Inner product, norm, and metric	9
64	2.4 Topology	11
65	2.4.1 Euclidean space	12
66	2.5 Measure theory	12
67	2.5.1 Sigma algebra	12
68	2.5.2 Measure	13
69	2.6 Integration	15
70	2.7 Functions	17
71	2.8 Pseudo-gradients	18
72	3 Game theory definitions	20
73	3.1 Normal-form game	20
74	3.2 Mixed-strategy game	22
75	3.3 Bayesian game	23
76	3.4 Mean-field game	24
77	3.5 Partially observable stochastic game	26
78	3.5.1 Episode formulation	27
79	3.5.2 Value formulation	27
80	4 Prior work	28
81	4.1 Prior work on infinite-action games	28
82	4.1.1 Action space discretization	28
83	4.1.2 Double oracle	29

84	4.1.3	Fictitious play	30
85	4.1.4	Evolution strategies	31
86	4.2	Prior work on infinite-player games	32
87	4.3	Prior work on learning dynamics	33
88	5	Completed work	38
89	5.1	Finding mixed-strategy equilibria of continuous-action games without gradients using randomized policy networks	39
90			
91	5.1.1	Method	39
92	5.1.2	Experiments	40
93	5.1.2.1	Continuous Colonel Blotto	41
94	5.1.2.2	Complete-information auction	43
95	5.1.2.3	Asymmetric-information auction	45
96	5.1.2.4	Chopstick auction	47
97	5.1.2.5	Visibility game	49
98	5.2	ApproxED: Approximate exploitability descent via learned best responses	52
99	5.2.1	Method	52
100	5.2.1.1	Best-response functions	52
101	5.2.1.2	Best-response ensembles	54
102	5.2.2	Experiments	55
103	5.2.2.1	Saddle-point game	56
104	5.2.2.2	Generalized matching pennies	56
105	5.2.2.3	Generalized rock paper scissors	57
106	5.2.2.4	Shapley game	57
107	5.2.2.5	Glicksberg–Gross game	58
108	5.2.2.6	Continuous security game	58
109	5.2.2.7	Poker games	60
110	5.2.2.8	GAN training	61
111	5.3	Joint-perturbation simultaneous pseudo-gradient	68
112	5.3.1	Method	68
113	5.3.2	Experiments	70
114	5.3.2.1	Multi-item unit-demand auction	70
115	5.3.2.2	Knapsack auction	72
116	5.3.2.3	Sequential auction for identical items	73
117	5.3.2.4	Continuous-action Goofspiel	73
118	5.4	Solving infinite-player games with player-to-strategy networks	75
119	5.4.1	Method	76
120	5.4.2	Experiments	79
121	5.4.2.1	Sum game	79
122	5.4.2.2	Anti-coordination game	79
123	5.4.2.3	Discontinuous game	81
124	5.4.2.4	Circle game	81
125	5.4.2.5	Cournot game	82
126	5.4.2.6	Bayesian Cournot game	84
127	5.4.2.7	Quadratic-cost Cournot game	86
128	5.4.2.8	Heterogeneous-cost Cournot game	87

129	5.4.2.9	Conformity game	88
130	5.4.2.10	Permutation game	88
131	6	Proposed work	91
132	6.1	Methods	91
133	6.1.1	Extension to multi-step games	91
134	6.1.2	Extension to infinite-step games	92
135	6.1.2.1	Continuous-time POSG formulation	92
136	6.1.2.2	Proposed approach	93
137	6.1.2.3	Open questions	93
138	6.1.2.4	Prior work	93
139	6.1.3	Using ApproxED-BRF to train the P2SN	94
140	6.1.4	Mitigating catastrophic forgetting in ApproxED-BRF	95
141	6.1.5	Modifying BRF’s update scheme to use other learning dynamics	96
142	6.1.6	Joint-perturbation estimator for other learning dynamics	97
143	6.2	Benchmarks	98
144	6.2.1	Epidemic: Epidemiological contagion and pandemic mitigation	99
145	6.2.2	Finance: Financial contagion and systemic risk	99
146	6.2.3	Grid: Wholesale forward capacity markets	100
147	6.2.4	Foraging: Evolutionary foraging with kin selection	101
148	6.2.5	Orbit: Orbital constellation management and debris avoidance	101
149	6.2.6	Traffic: Macroscopic urban traffic and fleet routing	102
150	6.3	Codebase	103
151	6.4	Available compute	103
152	6.5	Timeline	104
153	6.6	Stretch goals	105
154	6.6.1	Auction: Real-time programmatic ad auctions	106
155		Bibliography	107
156	A	Other completed work	127
157	B	Additional information	129
158	B.1	Additional information about auctions	129
159	B.2	Noise dimensionality in RPNs	130
160	B.3	Best response computation for continuous Colonel Blotto	131
161	C	Additional experiments	132
162	C.1	Ablation for Fourier features in P2SNs	132
163	C.2	Comparison to player discretization	135
164	D	Theorems	140
165	D.1	Equilibrium existence and uniqueness	140
166	D.1.1	Infinite-action games	140
167	D.1.2	Infinite-player games (existence)	141
168	D.1.3	Infinite-player games (uniqueness)	141
169	D.2	Exploitability	142

170	D.3 Subgradient descent	144
171	D.4 Theoretical results pertaining to P2SN	145
172	D.4.1 Progress guarantee	145
173	D.4.2 Special cases	146

174 Chapter 1

175 Introduction

176 1.1 Background

177 **Game theory** is the study of strategic interactions between agents. Roger Myerson, in *Game*
178 *Theory: Analysis of Conflict* (Myerson, 1991), defines it as “the study of mathematical models
179 of conflict and cooperation between intelligent rational decision-makers.” In these models, agents
180 may have independent interests, creating incentives for pure cooperation, pure competition, or a
181 complicated mixture of the two.

182 Within this framework, a **game** is a setting in which such agents interact. Each player has a
183 set of **strategies** to choose from. A **strategy profile** is a complete assignment of strategies to
184 players from their strategy sets. Finally, a **utility function** defines the total reward attained by
185 each player, given all players’ strategies.

186 The central solution concept in game theory is the **Nash equilibrium** (NE), named after
187 mathematician John Nash. It is a strategy profile for which no player has an incentive to unilaterally
188 change their strategy. In his foundational 1951 paper, *Non-Cooperative Games* (Nash, 1951), Nash
189 formally defined it as follows: “An equilibrium point is an n-tuple such that each player’s mixed
190 strategy maximizes his payoff if the strategies of the others are held fixed. Thus each player’s
191 strategy is optimal against those of the others.” Martin Osborne, in *An Introduction to Game Theory*
192 (Osborne, 2004), defines a Nash equilibrium as a strategy profile “with the property that no player
193 can do better by changing his strategy, given the other players’ strategies.”

194 **Machine learning** (ML) is a subset of artificial intelligence focused on building systems that
195 learn from data. In his foundational textbook *Machine Learning* (Mitchell, 1997), Tom Mitchell
196 provided the most widely accepted formal definition: “A computer program is said to learn from
197 experience E with respect to some class of tasks T and performance measure P , if its performance
198 at tasks in T , as measured by P , improves with experience E .” Most real-world problems have a
199 learnable structure, which is why inductive reasoning and ML work well (Goldblum et al., 2024).

200 **Neural networks** (NNs) are computational models inspired by biological neural architectures
201 and serve as a foundational pillar of modern ML. McCulloch and Pitts (1943) introduced the concept
202 of artificial neurons. Rosenblatt (1958) introduced the perceptron, an architecture and algorithm for
203 supervised learning of binary classifiers. A defining characteristic of neural networks is their capacity
204 to act as universal function approximators while maintaining a robust ability to generalize across
205 unseen inputs.

Specifically, the classical **universal approximation theorem** (UAT) (Cybenko, 1989; Hornik, 1991; Leshno et al., 1993; Pinkus, 1999) establishes that a feedforward network with a single, sufficiently wide hidden layer can approximate arbitrary continuous functions on a compact domain. While this guarantees the theoretical representational power of shallow networks, modern architectures typically rely on **depth** (multiple hidden layers) rather than extreme width to achieve highly efficient, hierarchical feature extraction.

Deep learning (DL) allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction (LeCun, Bengio, and Hinton, 2015). Subsequent extensions of the UAT specifically address DL, mathematically demonstrating that depth can exponentially reduce the required number of neurons for certain classes of functions (Telgarsky, 2016; Lin, Tegmark, and Rolnick, 2017). Furthermore, modern variants of the theorem establish that deep networks with bounded width (often strictly related to the input dimension) but arbitrary depth are also universal approximators (Lu et al., 2017; Yarotsky, 2017; Kidger and Lyons, 2020).

Reinforcement learning (RL) is a type of ML in which autonomous agents learn to make decisions by interacting with their environment. Richard Sutton and Andrew Barto, in their seminal work *Reinforcement Learning: An Introduction* (Sutton and Barto, 1998), defined RL as “learning what to do—how to map situations to actions—so as to maximize a numerical reward signal.” By receiving rewards for desired behaviors and penalties for undesirable ones, the agent discovers how to maximize cumulative rewards over time. It is widely used for complex, sequential tasks like game-playing, robotics, and personalized recommendations.

Multiagent reinforcement learning (MARL) sits at the intersection of game theory and ML. As described in *A Comprehensive Survey of Multiagent Reinforcement Learning* (Busoniu, Babuska, and Schutter, 2008), MARL investigates systems where “multiple agents learn by dynamically interacting with their environment and with each other.” Because the environment is made non-stationary by the concurrent learning of other agents, MARL algorithms must incorporate the strategic, game-theoretic reasoning required to anticipate and respond to other intelligent actors.

Real-world agents—humans and machines alike—are neither omniscient nor computationally unbounded. However, they are typically expected to act *approximately* rationally within their limitations. This idea, introduced by Herbert Simon in Simon (1955) and termed **bounded rationality**, holds that agents optimize subject to constraints on information, cognitive capacity, and time. This observation motivates solution concepts like ϵ -Nash equilibrium (Definition 85), in which each player’s regret is at most ϵ rather than exactly zero.

Despite human cognitive limitations, the rigorous mathematical prescriptions of game theory have proven highly effective when deployed by artificial agents against actual humans. For example, *DeepStack* (Moravčík et al., 2017) and *Libratus* (Brown and Sandholm, 2018) were the first AIs to achieve superhuman performance in heads-up (2-player) no-limit Texas Hold’em. Similarly, *Pluribus* (Brown and Sandholm, 2019) was the first AI to achieve superhuman performance in multi-player no-limit Texas Hold’em. By mastering the world’s most popular poker variant, these AIs demonstrated that game-theoretic algorithms can successfully navigate hidden information, deception, and human bounded rationality on a massive scale.

1.2 Goal

The main goal of this project is to develop a **unified algorithmic framework** capable of computing approximate Nash equilibria in highly complex, realistic game-theoretic environments. Specifically,

250 our goal is to design a framework that robustly handles games that have the following **complexities**:

251 • **Complex information structures:**

- 252 ○ Partial observability of the game state.
- 253 ○ Hidden information (players know things others do not).
- 254 ○ Imperfect recall (players might forget previously known information).

255 • **Infinite scale and dimensionality:**

- 256 ○ Infinitely many states.
- 257 ○ Infinitely many actions (and the computation of mixed strategies over them).
- 258 ○ Infinitely many players (in a general setting strictly beyond standard mean-field games, where
- 259 players are not necessarily symmetric or indistinguishable and may possess unique strategy
- 260 sets and payoffs).
- 261 ○ Infinitely many steps (e.g., differential games).

262 • **Intricate dynamics and payoffs:**

- 263 ○ Multiple steps (i.e., dynamic games).
- 264 ○ Stochastic observations, rewards, and state transitions.
- 265 ○ Discontinuous utility functions (e.g., to accurately model auctions).
- 266 ○ Matrix-valued discount factors (representing agents valuing each other’s future rewards,
- 267 naturally capturing social dynamics like altruism and spite).

268 Standard equilibrium-finding techniques fail to scale or even apply to this highly-complex setting.
269 We develop our framework to address all of these complexities simultaneously. We apply it to a suite
270 of **simulated environments** that mirror real-world applications. We evaluate it using standard
271 game-theoretic performance metrics. In particular, we estimate the empirical **exploitability** to
272 quantify the robustness of the computed strategies. To tackle the intractability of exact exploitability
273 computation in these domains, we estimate empirical exploitability as follows:

- 274 • Conditional sampling and action space discretization for single-step Bayesian games.
- 275 • Training best-response agents via reinforcement learning for multi-step dynamic games.

276 1.3 Motivation

277 In the next few subsections, we provide theoretical and practical motivations for some of the key
278 aspects of the very general setting we are tackling.

279 1.3.1 Infinite-action games

280 While most research on computing game-theoretic equilibria focuses on finite, discrete action spaces,
281 many real-world settings inherently involve continuous domains—such as allocations of **space**,
282 **time**, or **money**. Notable examples include continuous resource allocation games (Ganzfried, 2021),
283 security games in continuous spaces (Kamra, Fang, et al., 2017; Kamra, Gupta, Fang, et al., 2018;

284 Kamra, Gupta, Wang, et al., 2019), network games (Ghosh and Kundu, 2019), military simulations
285 and wargames (Marchesi, Trovò, and Gatti, 2020), and complex video games (Berner et al., 2019;
286 Vinyals et al., 2019). Furthermore, even when an action space is technically discrete, it can be so
287 fine-grained that treating it as continuous yields significant computational efficiency (Borel and
288 Ville, 1938; Chen and Ankenman, 2006; Ganzfried and Sandholm, 2010).

289 These infinite-action games often introduce additional complexities. For example, **discontinuous**
290 **utility functions** are highly relevant to real-world scenarios, particularly in economic settings
291 like auctions. Finally, because players must often randomize their actions to avoid being exploited
292 (assuming access to randomization oracles), it is essential to be able to represent and learn **mixed**
293 **strategies** over these infinite action spaces.

294 1.3.2 Infinite-player games

295 Many disciplines—ranging from economics and financial markets to congestion analysis, crowd
296 dynamics, epidemiology, and population ecology—study systems comprising **large numbers of**
297 **interacting agents**. The presence of numerous entities, each pursuing individual and often competing
298 interests, makes modeling and analyzing these multiagent systems notoriously complex. To simplify
299 this analysis, researchers frequently model such systems as *infinite*-player games. Taking the limit as
300 the number of players approaches infinity yields a mathematical tractability akin to the treatment
301 of **many-particle systems** in thermodynamics and statistical mechanics.

302 This continuous approximation has a rich history. A foundational paper by Aumann (1964)
303 proposed modeling markets with individually insignificant traders as a continuum, analogous to the
304 continuous points on a line. Building on this concept, **mean-field game** (MFG) theory studies
305 the limit of symmetric N -player games as N approaches infinity (Lasry and Lions, 2007; Huang,
306 Malhamé, and Caines, 2006; Huang, Caines, and Malhamé, 2007). In an MFG, agents are individually
307 negligible and exchangeable. Rather than tracking every pairwise interaction, each infinitesimal
308 agent optimizes its control against a **mean field** representing the aggregate state distribution of the
309 population. Equilibrium is thus defined as a self-consistent mean field generated by the agents’ best
310 responses.

311 Computing Nash equilibria of these many-agent systems is invaluable not only for an individual
312 agent determining its optimal strategy, but also for social scientists seeking to understand the system
313 from the outside, and for social planners aiming to direct it. In these contexts, the **aggregate**
314 **behavior** of players is often what truly matters, rather than the specific actions of individuals. In
315 mechanism design, for instance, social planners strive to engineer games whose equilibria possess
316 desired properties.

317 Consequently, this is highly applicable to real-world settings such as many-player auctions,
318 financial markets, and transit networks. From a theoretical perspective, it can also be utilized to
319 study the equilibria of environments like **abstract economies**. Ultimately, this framework can be
320 employed both to make precise prescriptions for how individual agents should act rationally, *and* to
321 predict the emergent behavior of self-interested agents within a massive system.

322 1.3.3 Matrix-valued discount factors

323 A **discount factor** quantifies how much an agent values future rewards relative to immediate
324 ones. In a multiagent game, each player can have a distinct discount factor. Traditionally, these are
325 represented as a *discount vector*, where each component corresponds to an individual player.

326 This concept can be generalized to *matrix-valued* discount factors. The diagonal elements of such
 327 a matrix represent the standard, individual discount factors, while the off-diagonal elements capture
 328 cross-player interactions. A non-zero off-diagonal factor indicates that an agent internalizes *another*
 329 agent’s future rewards, either positively or negatively.

330 Assuming a parameter $\alpha > 0$, we can model a variety of distinct behavioral dynamics between
 331 two agents using the following 2×2 matrices.

- 332 • **Egoism (strict self-interest)**: Both agents are strictly self-interested and perfectly indifferent
 333 to their counterpart’s outcomes.

$$\gamma_{\text{egoists}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (1.1)$$

- 334 • **Lovers (mutual altruism)**: Both agents value their counterpart’s happiness in addition to
 335 their own.

$$\gamma_{\text{lovers}} = \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix} \quad (1.2)$$

- 336 • **Bitter enemies (mutual spite)**: Both agents actively derive disutility from their counterpart’s
 337 success.

$$\gamma_{\text{enemies}} = \begin{bmatrix} 1 & -\alpha \\ -\alpha & 1 \end{bmatrix} \quad (1.3)$$

- 338 • **The unrequited lover (one-sided altruism)**: Player 1 values Player 2’s happiness, but Player
 339 2 is entirely indifferent to Player 1.

$$\gamma_{\text{unrequited}} = \begin{bmatrix} 1 & \alpha \\ 0 & 1 \end{bmatrix} \quad (1.4)$$

- 340 • **The saboteur (one-sided spite)**: Player 1 actively dislikes Player 2 and seeks to minimize
 341 their counterpart’s reward, while Player 2 remains a pure egoist.

$$\gamma_{\text{saboteur}} = \begin{bmatrix} 1 & -\alpha \\ 0 & 1 \end{bmatrix} \quad (1.5)$$

- 342 • **Hyper-competitive (zero-sum mindset)**: The agents care only about the *relative difference*
 343 in their rewards. A counterpart’s gain is mathematically identical to a personal loss.

$$\gamma_{\text{zero-sum}} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \quad (1.6)$$

- 344 • **Martyrs (pure self-sacrifice, self-abnegating altruism)**: Both agents place zero value on
 345 their own rewards and care *only* about the outcomes of their counterpart.

$$\gamma_{\text{martyrs}} = \begin{bmatrix} 0 & \alpha \\ \alpha & 0 \end{bmatrix} \quad (1.7)$$

346 In game theory and sociobiology, these interactive preferences are commonly studied under the
347 name of **altruism and spite**. There is a rich body of literature that explores how such motives
348 cause behavior to deviate from strict self-interest.

349 In behavioral and experimental economics, this class of preferences is known as other-regarding or
350 **social preferences**. Fehr and Schmidt (1999) proposed the canonical model of **inequity aversion**,
351 in which an agent’s utility is penalized by both advantageous and disadvantageous deviations from
352 their counterpart’s payoff. Bolton and Ockenfels (2000) proposed the related ERC (equity, reciprocity,
353 and competition) framework, formulated in terms of relative rather than absolute comparisons.
354 Charness and Rabin (2002) tested these models experimentally, disentangling inequity aversion
355 from broader concerns for social welfare and reciprocity. Collectively, these models motivate the
356 explicit parameterization of interpersonal weightings in utility. This is precisely the role played by
357 the off-diagonal entries of our matrix-valued discount factors.

358 For instance, in the context of public goods, Saijo and Nakamura (1995) examined deviations
359 from strict economic rationality—where pure self-interest dictates either full contribution or free-
360 riding—highlighting the role of spite in voluntary contribution mechanisms. Levine (1998) proposed
361 a broader theory of altruism where a player’s payoff is a linear combination of their own and their
362 counterpart’s monetary income. By analyzing ultimatum and centipede game experiments, they
363 estimated the population distribution of altruism and spite, noting that the weight placed on a
364 counterpart’s income relies on private information and beliefs about the counterpart’s disposition.

365 The impact of spite is particularly well-documented in **auction theory**, where agents frequently
366 bid more aggressively than standard models prescribe. This aggressiveness can stem from inherent
367 psychological spite or strategic advantages, such as long-term benefits in closed markets where a
368 competitor’s loss translates to personal gain. Morgan, Steiglitz, and Reis (2003) demonstrated that
369 bidders who attach disutility to rival surplus bid aggressively enough to rationalize anomalies observed
370 in laboratory settings. Brandt, Sandholm, and Shoham (2007) corroborated this by analyzing spiteful
371 behavior in sealed-bid auctions. Expanding on these foundations, Sharma and Sandholm (2010)
372 introduced the first analysis of asymmetric spite, reflecting the realistic scenario where bidders are
373 spiteful to varying degrees. Finally, addressing the ultimate mechanism design challenge in this
374 domain, Tang and Sandholm (2012) derived optimal auctions tailored specifically for both spiteful
375 and altruistic bidders.

376 Beyond static or single-step interactions, the concept of generalized discounting has been extended
377 to sequential decision-making and dynamic environments. In the context of temporal discounting,
378 White (2017) introduced a reinforcement learning task formalism that generalizes the traditional
379 scalar discount factor into a **transition-dependent variable**. Building on this flexibility, Pitis
380 (2019) provided an axiomatic, decision-theoretic proof demonstrating that standard, fixed scalar
381 discounting is fundamentally insufficient for modeling general human preferences. This expanded
382 framework accommodates a state-action-dependent discount factor that **can be greater than 1**,
383 provided there is eventual long-run discounting.

384 Parallel advancements have formalized social discounting in multi-step settings, where the
385 matrix values represent **interpersonal weights** rather than temporal horizons. McKee et al. (2020)
386 integrated **social value orientation** (SVO) directly into deep multiagent reinforcement learning.
387 They demonstrated that introducing heterogeneity in SVOs generates meaningful and complex
388 behavioral variation among agents, closely mirroring the predictions of interdependence theory.
389 Formalizing the mathematical limits of these interactions, Ray and Vohra (2020) studied strategic
390 environments characterized by **payoff-based externalities**—situations where an agent’s utility
391 depends not only on their own actions but directly on the realized payoffs of others. They termed

392 the resulting systems “**games of love and hate**”, providing a rigorous foundation for the infinitely
393 reflecting utilities that arise within these interdependent matrices.

394 1.4 Outline

395 The structure of this thesis proposal is as follows. In Section 2, we define the mathematical notation
396 and concepts that are required for our setting, which deals with infinite spaces. In Section 3, we
397 describe the game-theoretic concepts that are needed for our setting, and give a mathematically
398 rigorous formulation of the problem we are tackling. In Section 4, we describe prior work that has
399 been done in each of the main areas we are tackling. In Section 5, we describe the work we have
400 already completed that contributes toward our ultimate goal. In Section 6, we describe the next
401 steps for our project and proposed future work. In the appendix, we include additional material.
402 This includes additional experiments and theoretical analysis. Section A describes other works I
403 completed during my PhD that fall outside the scope of this thesis.

404 **Chapter 2**

405 **Mathematical preliminaries**

406 In this section, we define the mathematical notation and concepts that are required for our setting.
 407 Since we deal with infinite spaces, we need to define related concepts rigorously.

408 **2.1 Notation**

409 In the expressions below, A and B are sets, f and g are functions, n is a natural number, ϕ is a
 410 formula, X is a random variable, d is a metric, and V is a Banach space.

$x \in A$	x is an element of A	
$A \subseteq B$	A is a subset of B	$\forall x \in A, x \in B$
$A \cup B$	union of A and B	$\{x \mid x \in A \vee x \in B\}$
$A \cap B$	intersection of A and B	$\{x \mid x \in A \wedge x \in B\}$
$A \setminus B$	relative complement of B in A	$\{x \mid x \in A \wedge x \notin B\}$
$\mathcal{P}(A)$	powerset of A	$\{B \mid B \subseteq A\}$
$\cup A$	union of A	$\{x \mid \exists B \in A, x \in B\}$
\emptyset	empty set	$\{x \mid \perp\}$
$\sup A$	supremum of A	
$\inf A$	infimum of A	
$\text{dom } f$	domain of f	
$\text{cod } f$	codomain of f	
$A \times B$	Cartesian product of A and B	
$A \rightarrow B$ or B^A	set of functions from A to B	
$f[A]$	image of A under f	$\{f(x) \mid x \in A\}$
$f^{-1}[A]$	preimage of A under f	$\{x \mid f(x) \in A\}$
$\text{im } f$	image of f	$f[\text{dom } f]$
$\text{zero } f$	zeroset of f	$f^{-1}[\{0\}]$
$\text{supp } f$	support of f	$(\text{dom } f) \setminus (\text{zero } f)$
$\text{argmax } f$	argmax of f	$f^{-1}[\{\sup(\text{im } f)\}]$
$\text{argmin } f$	argmin of f	$f^{-1}[\{\inf(\text{im } f)\}]$
Πf	Cartesian product of f	
$f \circ g$	composition of f and g	

\mathbb{N}	set of natural numbers	
\mathbb{Z}	set of integers	
\mathbb{Q}	set of rational numbers	
\mathbb{R}	set of real numbers	
$\overline{\mathbb{R}}$	set of extended real numbers	$\mathbb{R} \cup \{-\infty, \infty\}$
$[n]$	set of natural numbers less than n	$\{0, \dots, n-1\}$
$\llbracket \phi \rrbracket$	Iverson bracket of ϕ	
$\mathbb{E} X$	expected value of X	
$\text{cl } A$	topological closure of A	
$\text{int } A$	topological interior of A	
∂A	topological boundary of A	$(\text{cl } A) \setminus (\text{int } A)$
$\mathbb{B}_r(x)$	open ball of radius r centered at x	$\{y \mid d(x, y) < r\}$
$\mathcal{L}(V)$	space of bounded linear operators on V	

411 2.2 Set theory

412 **Definition 1.** The **Cartesian product** of an indexed family of sets is defined by

$$\prod f = \{g : \text{dom } f \rightarrow \cup \text{im } f \mid \forall x \in \text{dom } g, g(x) \in f(x)\} \quad (2.1)$$

413 **Definition 2.** Let $\{X_i\}_{i \in I}$ be an indexed family of sets. For any $i \in I$, the **projection map** onto i
414 is the map $\pi_i : \prod_{i \in I} X_i \rightarrow X_i$ such that $\pi_i(x) = x_i$.

415 **Definition 3.** Let $\{X_i\}_{i \in I}$ be an indexed family of sets. For any $J \subseteq I$, the **canonical projection**
416 **map** onto J is the map $\pi_J : \prod_{i \in I} X_i \rightarrow \prod_{j \in J} X_j$ such that $\pi_J(\{x_i\}_{i \in I}) = \{x_j\}_{j \in J}$.

417 **Definition 4.** Let $\{x_i\}_{i \in I}$ be an indexed family. Then x_{-i} denotes $x|_{I \setminus \{i\}}$.

418 **Definition 5.** A **set family** is a tuple (X, F) in which X is a set and $F \subseteq \mathcal{P}(X)$. That is,
419 $\forall A \in F, A \subseteq X$.

420 **Definition 6.** A **ring of sets** is a set family closed under binary union and relative complement.
421 That is, R is a ring of sets iff $\forall A, B \in R, A \cup B \in R \wedge A \setminus B \in R$.

422 **Definition 7.** If A is a set, the **indicator** function 1_A is the function defined by

$$1_A(x) = \llbracket x \in A \rrbracket \quad (2.2)$$

423 **Definition 8.** A **preorder** is a binary relation that is reflexive and transitive.

424 **Definition 9.** Let \preceq be a preorder on a set P . Let $S \subseteq P$ and $s \in P$. s is a **supremum** or **least**
425 **upper bound** of S iff $\forall u \in P (s \preceq u \leftrightarrow \forall x \in S, x \preceq u)$. s is an **infimum** or **greatest lower bound**
426 of S iff it is a supremum under the reverse preorder \succeq .

427 2.3 Inner product, norm, and metric

428 Let K be a subring of \mathbb{C} that is closed under complex conjugation. Let V be a module over K .

429 **Definition 10.** An **inner product** is a function $\langle \cdot, \cdot \rangle : V \times V \rightarrow K$ that satisfies the following
430 properties for all $x, y, z \in V$ and $\alpha \in K$:

- 431 1. Conjugate symmetry: $\langle x, y \rangle = \langle y, x \rangle^*$.
- 432 2. Linearity: $\langle \alpha x + y, z \rangle = \alpha \langle x, z \rangle + \langle y, z \rangle$.
- 433 3. Positive-definiteness: $\langle x, x \rangle \geq 0$, with equality iff $x = 0$.

434 **Proposition 1.** The **Cauchy–Schwarz inequality** (Cauchy, 1821b; Bunyakovsky, 1859; Schwarz,
435 1890) states that for all $x, y \in V$:

$$|\langle x, y \rangle|^2 \leq \langle x, x \rangle \langle y, y \rangle \quad (2.3)$$

436 Here, $|\cdot|$ denotes the complex modulus on \mathbb{C} .

437 **Definition 11.** A **norm** is a function $\|\cdot\| : V \rightarrow \mathbb{R}$ that satisfies the following properties for all
438 $x, y \in V$ and $\alpha \in K$:

- 439 1. Positive-definiteness: $\|x\| \geq 0$, with equality iff $x = 0$.
- 440 2. Absolute homogeneity: $\|\alpha x\| = |\alpha| \|x\|$.
- 441 3. Triangle inequality: $\|x + y\| \leq \|x\| + \|y\|$.

442 **Definition 12.** The **induced norm** of an inner product $\langle \cdot, \cdot \rangle$ is $\|x\| = \sqrt{\langle x, x \rangle}$. The triangle
443 inequality follows from the Cauchy–Schwarz inequality.

444 **Definition 13.** Let X be a set. A **metric** on X is a function $d : X \times X \rightarrow \mathbb{R}$ that satisfies the
445 following properties for all $x, y, z \in X$:

- 446 1. Positive-definiteness: $d(x, y) \geq 0$, with equality iff $x = y$.
- 447 2. Symmetry: $d(x, y) = d(y, x)$.
- 448 3. Triangle inequality: $d(x, z) \leq d(x, y) + d(y, z)$.

449 **Definition 14.** The **induced metric** of a norm $\|\cdot\|$ is $d(x, y) = \|x - y\|$.

450 **Definition 15.** A **Cauchy sequence** (Cauchy, 1821a) is a sequence $\{x_n\}_{n \in \mathbb{N}}$ such that

$$\forall \varepsilon > 0, \exists N \in \mathbb{N}, \forall i, j > N, d(x_i, x_j) < \varepsilon \quad (2.4)$$

451 **Definition 16.** A **complete space** is a space in which every Cauchy sequence has a limit.

452 **Definition 17.** A **Banach space** is a complete normed vector space.

453 2.4 Topology

454 **Definition 18.** A **topological space** is a set family (X, τ) where τ is closed under finite intersections
 455 and arbitrary unions. Elements of τ are called **open sets**.

456 **Definition 19.** Let (X, τ_X) and (Y, τ_Y) be topological spaces. A function $f : X \rightarrow Y$ is **continuous**
 457 iff the preimage of every open set is an open set. That is,

$$\forall U \in \tau_Y, f^{-1}[U] \in \tau_X \quad (2.5)$$

458 **Definition 20.** Let $\{(X_i, \tau_i)\}_{i \in I}$ be a collection of topological spaces. The **product topology** is
 459 the smallest topology on $\prod X$ such that every projection map is continuous.

460 **Proposition 2.** Under the product topology, for any topological space Y , a function $f : Y \rightarrow \prod_{i \in I} X_i$
 461 is continuous iff, for each $i \in I$, the component function $\pi_i \circ f$ is continuous.

462 **Definition 21.** Let (X, τ) be a topological space. A **base** (or **basis**) for τ is a family $\mathcal{B} \subseteq \tau$ of open
 463 sets such that every open set in the topology can be represented as the union of some subfamily of
 464 \mathcal{B} :

$$(\forall U \in \tau)(\exists \mathcal{F} \subseteq \mathcal{B})(\cup \mathcal{F} = U) \quad (2.6)$$

465 **Definition 22.** Let X be a set. A family \mathcal{B} of subsets of X satisfies the **base axioms** (and thus
 466 forms a valid base for a topology on X) iff it satisfies the following two conditions:

467 1. **Covering:** The union of all sets in \mathcal{B} is the entire set X .

$$\bigcup \mathcal{B} = X \quad (2.7)$$

468 Equivalently, for every $x \in X$, there exists at least one $B \in \mathcal{B}$ such that $x \in B$.

469 2. **Intersection:** If $B_1, B_2 \in \mathcal{B}$ and x is a point in their intersection ($x \in B_1 \cap B_2$), then there
 470 exists a third basis set $B_3 \in \mathcal{B}$ that contains x and fits entirely inside the intersection:

$$\exists B_3 \in \mathcal{B} : x \in B_3 \subseteq B_1 \cap B_2 \quad (2.8)$$

471 **Definition 23.** Let X be a set and let \mathcal{B} be a family of subsets of X that satisfies the base axioms.
 472 The **topology generated** by \mathcal{B} is the collection of all unions of subfamilies of \mathcal{B} :

$$\tau = \left\{ \bigcup \mathcal{F} \mid \mathcal{F} \subseteq \mathcal{B} \right\} \quad (2.9)$$

473 **Definition 24.** A **topological field** is a field equipped with a topology such that all field operations
 474 (addition, multiplication, and inversion on non-zero elements) are continuous.

475 **Definition 25.** A **topological vector space** (TVS) is a vector space over a topological field
 476 equipped with a topology such that the vector addition and scalar multiplication maps are continuous.

477 **Definition 26.** Let (X, τ) be a topological space. Let $x \in X$ be a point. An **open neighborhood**
 478 of x is an open set $U \in \tau$ such that $x \in U$.

479 **Definition 27.** A **locally convex topological vector space** (LCTVS) is a TVS in which every
 480 open neighborhood of the zero vector contains a convex open neighborhood of the zero vector. That
 481 is, it is a TVS that possesses a local base at the zero vector consisting entirely of convex sets.

482 **Definition 28.** Two points $x, y \in X$ can be **separated** by open neighborhoods iff there exist open
 483 sets $U, V \in \tau$ that are disjoint ($U \cap V = \emptyset$) and $x \in U, y \in V$.

484 **Definition 29.** A topological space **separates points** (equivalently, is a **Hausdorff space** or T_2
 485 space) iff any two distinct points can be separated by open neighborhoods.

486 2.4.1 Euclidean space

487 **Definition 30.** The **Euclidean inner product** is

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i \quad (2.10)$$

488 **Definition 31.** The **Euclidean norm** is the norm induced by this inner product.

489 **Definition 32.** The **Euclidean metric** is the metric induced by this norm.

490 **Definition 33.** The **Euclidean basis** is the collection of open balls induced by this metric:

$$B = \{\mathbb{B}_r(x) \mid x \in \mathbb{R}^n \wedge r \in \mathbb{R} \wedge r > 0\} \quad (2.11)$$

491 **Definition 34.** The **Euclidean topology** is the topology generated by this basis.

492 2.5 Measure theory

493 2.5.1 Sigma algebra

494 **Definition 35.** Let X be a set. A **σ -algebra** is a collection $\Sigma \subseteq \mathcal{P}(X)$ where $X \in \Sigma$ and Σ is
 495 closed under complements and countable unions.

496 **Definition 36.** A **measurable space** is a pair (X, Σ) where Σ is a σ -algebra on X . Elements of Σ
 497 are called **measurable sets**.

498 **Definition 37.** Let (X, Σ_X) and (Y, Σ_Y) be measurable spaces. A function $f : X \rightarrow Y$ is **measur-**
 499 **able** iff the preimage of every measurable set is a measurable set. That is,

$$\forall E \in \Sigma_Y, f^{-1}[E] \in \Sigma_X \quad (2.12)$$

500 **Definition 38.** Let $\{(X_i, \Sigma_i)\}_{i \in I}$ be a collection of measurable spaces. The **product σ -algebra**
 501 is the smallest σ -algebra on $\prod X$ such that every projection map is measurable. It is denoted by
 502 $\bigotimes_{i \in I} \Sigma_i$.

503 **Proposition 3.** Under the product σ -algebra, for any measurable space Y , a function $f : Y \rightarrow$
 504 $\prod_{i \in I} X_i$ is measurable iff, for each $i \in I$, the component function $\pi_i \circ f$ is measurable.

505 **Definition 39.** Let X be a set and (Y, Σ) be a measurable space. For each $x \in X$, the evaluation
 506 functional $e_x \in Y^X \rightarrow Y$ is defined by $e_x(f) = f(x)$. The **σ -algebra generated by the evaluation**
 507 **functionals** is the product σ -algebra on Y^X .

508 **Definition 40.** Let (X, F) be a set family. The σ -algebra generated by F is the smallest σ -algebra
 509 that includes F :

$$\sigma(F) = \cap\{\Sigma \subseteq \mathcal{P}(X) \mid F \subseteq \Sigma \wedge \Sigma \text{ is a } \sigma\text{-algebra}\} \quad (2.13)$$

510 **Definition 41.** Let (X, τ) be a topological space. The **Borel σ -algebra** is the σ -algebra generated
 511 by τ : $\mathcal{B}(X) = \sigma(\tau)$. Elements of $\mathcal{B}(X)$ are called **Borel sets**.

512 To equip sets with σ -algebras, we adopt the following conventions.

- 513 • **Topological spaces:** If X is a topological space, it is equipped with the Borel σ -algebra $\mathcal{B}(X)$.
- 514 • **Banach spaces:** If V is a Banach space, it is equipped with the Borel σ -algebra generated by
 515 its norm topology.
- 516 • **Product spaces:** If $\{(X_i, \Sigma_i)\}_{i \in I}$ is a collection of measurable spaces, the Cartesian product
 517 $\prod_{i \in I} X_i$ is equipped with the product σ -algebra.
- 518 • **Sets of measures:** If (X, Σ) is a measurable space, any set of measures is equipped with the
 519 σ -algebra generated by the evaluation functionals.

520 2.5.2 Measure

521 **Definition 42.** Let R be a ring of sets. A **pre-measure** on R is a function $\nu : R \rightarrow [0, \infty]$ such
 522 that $\nu(\emptyset) = 0$ and ν is countably additive.

523 **Definition 43.** A **measure** is a **pre-measure** on a σ -algebra.

524 **Definition 44.** A **measure space** is a triple (X, Σ, μ) where μ is a measure on Σ .

525 **Definition 45.** A **probability measure** is a measure μ where $\mu(X) = 1$.

526 **Definition 46.** A **probability space** is a measure space with a probability measure.

527 **Definition 47.** A **Borel probability measure** on X is a probability measure on $(X, \mathcal{B}(X))$.

528 **Definition 48.** The set of Borel probability measures on X is denoted by ΔX .

529 **Definition 49.** The **Dirac measure** at x is denoted by δ_x . It is defined by $\delta_x(A) = \llbracket x \in A \rrbracket$.

530 **Definition 50.** Let (X, Σ_X, μ) be a measure space, (Y, Σ_Y) be a measurable space, and $f : X \rightarrow Y$
 531 be a measurable function. The **pushforward measure** of μ by f , denoted by $f_*\mu$, is the measure
 532 on (Y, Σ_Y) defined by

$$\forall B \in \Sigma_Y, f_*\mu(B) = \mu(f^{-1}[B]) \quad (2.14)$$

533 **Definition 51.** The **Carathéodory extension** of a pre-measure ν is the restriction of the outer
 534 measure μ^* of ν to the σ -algebra of μ^* -measurable sets.

535 **Definition 52.** Let $\{(X_i, \Sigma_i, \mu_i)\}_{i \in I}$ be a collection of measure spaces. The **product measure**
 536 $\bigotimes_{i \in I} \mu_i$ is the Carathéodory extension of the set function ν such that for all $\{E_i\}_{i \in I}$ where $E_i \in \Sigma_i$
 537 for all i and $E_i = X_i$ for all but finitely many i ,

$$\nu\left(\prod_{i \in I} E_i\right) = \prod_{i \in I} \mu_i(E_i) \quad (2.15)$$

538 **Definition 53.** Let $\{(X_i, \Sigma_i, \mu_i)\}_{i \in I}$ be a countable collection of measure spaces. The **sum measure**
 539 $\sum_{i \in I} \mu_i$ is the measure ν on the disjoint union $\bigsqcup_{i \in I} X_i$ such that

$$\nu(E) = \sum_{i \in I} \mu_i(E \cap X_i) \quad (2.16)$$

540 **Definition 54.** Let (X, Σ, μ) be a measure space and $A \in \Sigma$ be a measurable set. The **restriction**
 541 **measure** $\mu|_A \in \Sigma \rightarrow [0, \infty]$ is a new measure defined by

$$\forall B \in \Sigma, \mu|_A(B) = \mu(A \cap B) \quad (2.17)$$

542 **Definition 55.** Let $\{(X_i, \Sigma_i, \mu_i)\}_{i \in I}$ be a collection of measure spaces. A **coupling** is a measure γ
 543 on the product measurable space $(\prod_{i \in I} X_i, \bigotimes_{i \in I} \Sigma_i)$ such that the **marginal** measure along every
 544 i -th coordinate space is μ_i :

$$\forall i \in I, (\pi_i)_* \gamma = \mu_i \quad (2.18)$$

545 where $\pi_i \in \prod_{j \in I} X_j \rightarrow X_i$ is the i -th projection map.

546 We also have the following definitions.

A null	$A \in \Sigma \wedge \mu(A) = 0$
B negligible	$\exists A, B \subseteq A \wedge A$ null
A co-null	$X \setminus A$ null
B co-negligible	$X \setminus B$ negligible
547 ess sup f	$\inf\{a \in \mathbb{R} \mid \{x \in X \mid f(x) > a\} \text{ negligible}\}$
ess inf f	$\sup\{a \in \mathbb{R} \mid \{x \in X \mid f(x) < a\} \text{ negligible}\}$
ϕ almost nowhere	$\{x \in X \mid \phi(x)\} \text{ negligible}$
ϕ almost everywhere	$\{x \in X \mid \neg\phi(x)\} \text{ negligible}$

548 **Definition 56.** An **atom** is a measurable set $A \in \Sigma$ such that $\mu(A) > 0$ and, for every measurable
 549 subset $B \subseteq A$, $\mu(B) = 0$ or $\mu(B) = \mu(A)$.

550 **Definition 57.** A measure is **atomic** iff every measurable set of positive measure contains an atom.

551 **Definition 58.** A measure is **non-atomic**, **atomless**, or **diffuse** iff it has no atoms.

552 **Definition 59.** Let (M, d) be a metric space. Let μ and ν be probability measures on M . The
 553 **Wasserstein distance** of order p between μ and ν is

$$W_p(\mu, \nu) = \left(\inf_{\gamma \in \Gamma(\mu, \nu)} \mathbb{E}_{(x, y) \sim \gamma} d(x, y)^p \right)^{\frac{1}{p}} \quad (2.19)$$

554 where $\Gamma(\mu, \nu)$ is the set of all couplings of μ and ν .

555 **Proposition 4.** The case $p \rightarrow \infty$ is as follows.

$$W_\infty(\mu, \nu) = \lim_{p \rightarrow \infty} W_p(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \operatorname{ess\,sup}_{(x,y) \sim \gamma} d(x, y) \quad (2.20)$$

556 We adopt the following canonical conventions for equipping new spaces with measures:

- 557 • **Products:** If each X_i has a measure μ_i , we equip $\prod_{i \in I} X_i$ with the product measure $\bigotimes_{i \in I} \mu_i$.
- 558 • **Mappings:** If f is a measurable function and $\operatorname{dom} f$ has a measure μ , we equip $\operatorname{cod} f$ with the
559 pushforward measure $f_*\mu$.
- 560 • **Subsets:** If X has a measure μ and $A \subseteq X$ is a measurable subset, we equip A with the
561 restriction measure $\mu|_A$.
- 562 • **Disjoint unions:** If I is countable and each X_i has a measure μ_i , we equip $\bigsqcup_{i \in I} X_i$ with the
563 sum measure $\sum_{i \in I} \mu_i$.
- 564 • **Quotients:** If \sim is an equivalence relation on X and X has a measure μ , we equip the quotient
565 space X/\sim with the pushforward measure $\pi_*\mu$, where $\pi : X \rightarrow X/\sim$ is the quotient map.

566 2.6 Integration

567 **Definition 60.** The **integral** of a function f with respect to a measure μ is denoted by $\int_{x \sim \mu} f(x)$
568 or $\int_\mu f$. If μ is a probability measure, we also use the notation for expectations, $\mathbb{E}_{x \sim \mu} f(x)$ or $\mathbb{E}_\mu f$.

569 An integral satisfies the following fundamental properties:

570 **Proposition 5** (Indicator functions). Let A be a measurable set. Then

$$\int_\mu 1_A = \mu(A) \quad (2.21)$$

571 **Proposition 6** (Addition). Let f and g be integrable functions. Then

$$\int_\mu (f + g) = \int_\mu f + \int_\mu g \quad (2.22)$$

572 **Proposition 7** (Scalar multiplication). Let f be an integrable function and α be a scalar. Then

$$\int_\mu \alpha f = \alpha \int_\mu f \quad (2.23)$$

573 **Definition 61.** A **simple function** is a function of the form

$$f = \sum_{i=1}^n \alpha_i 1_{A_i} \quad (2.24)$$

574 where $n \in \mathbb{N}$, each α_i is a scalar, and each A_i is a measurable set.

575 It follows from the preceding properties that any simple function, provided its constituent sets
576 A_i have finite measure, is integrable.

577 There are many types of integrals in the literature, including (but not limited to):

- 578 • The Riemann integral (Riemann, 1868)
- 579 • The Lebesgue integral (Lebesgue, 1902)
- 580 • The Riemann–Stieltjes integral (Stieltjes, 1894)
- 581 • The Lebesgue–Stieltjes integral (Saks, 1937)
- 582 • The Henstock–Kurzweil integral (Kurzweil, 1957)
- 583 • The Bochner integral (Bochner, 1933)
- 584 • The Pettis integral (Pettis, 1938)
- 585 • The Dunford integral (Dunford, 1935)

586 Some of these integrals are **stronger** than others; that is, they can integrate a **strictly larger**
 587 class of functions. For example, in addition to linearity and the indicator property, the **Bochner**
 588 **integral** satisfies the following strong continuity guarantee:

589 **Proposition 8** (Continuity). Let (Ω, Σ, μ) be a measure space. Let B be a Banach space with norm
 590 $\|\cdot\|$. Let $f : \Omega \rightarrow B$. For each $n \in \mathbb{N}$, let $g_n : \Omega \rightarrow B$ be Bochner-integrable functions. Then

$$\lim_{n \rightarrow \infty} \int_{x \sim \mu} \|g_n(x) - f(x)\| = 0 \implies \lim_{n \rightarrow \infty} \int_{\mu} g_n = \int_{\mu} f \quad (2.25)$$

591 This means that if a sequence of functions converges to f in the mean (i.e., the integral of their
 592 normed difference vanishes), then their respective integrals converge to the integral of f .

593 An even more generalized formulation is the **Pettis integral**, which is entirely characterized by
 594 the following property:

595 **Proposition 9** (Duality). Let (Ω, Σ, μ) be a measure space. Let V be a LCTVS with continuous
 596 dual space V^* . Let $f : \Omega \rightarrow V$ be a function. Let $v \in V$. Then

$$\left(\forall \phi \in V^* : \phi(v) = \int_{\mu} (\phi \circ f) \right) \implies \int_{\mu} f = v \quad (2.26)$$

597 In functional analysis, this leads directly to the core property that characterizes all such vector-
 598 valued integrals.

599 **Definition 62.** The **universal property of integration** is defined as follows. Let (Ω, Σ, μ) be a
 600 measure space. Let X and Y be LCTVSs over the same scalar field. If $f : \Omega \rightarrow X$ is an integrable
 601 function, its integral is the unique vector $\int_{\mu} f \in X$ such that, for any continuous linear operator
 602 $\phi : X \rightarrow Y$, the composition $\phi \circ f : \Omega \rightarrow Y$ is also integrable, and

$$\phi \left(\int_{\mu} f \right) = \int_{\mu} (\phi \circ f) \quad (2.27)$$

603 **Proposition 10.** The universal property completely determines the integral. Evaluating it over the
 604 continuous dual space (by letting Y be the underlying scalar field, so that $\phi \in X^*$) uniquely isolates
 605 the vector in X , provided the topology of X separates points.

2.7 Functions

Let f be a function.

Definition 63. Let S be a set. The **restriction** of f to S is denoted by $f|_S$. It is equal to $f \cap (S \times \text{im } f)$, by the standard definition of a function as a set of ordered input-output pairs.

Definition 64. The **substitution** $f[a \mapsto b]$ is the function defined by

$$f[a \mapsto b](x) = \begin{cases} b & x = a \\ f(x) & x \neq a \end{cases} \quad (2.28)$$

The upper-right **Dini derivative** of f at x in direction d is

$$(D^+ f)(x)(d) = \limsup_{t \downarrow 0} \frac{f(x + td) - f(x)}{t} \quad (2.29)$$

If f is differentiable, this is just $\langle \nabla f(x), d \rangle$.

Definition 66. f is **convex** iff for all $x, y \in \text{dom } f$ and all $t \in [0, 1]$:

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) \quad (2.30)$$

Definition 67. f is **m -convex** iff $x \mapsto f(x) - \frac{m}{2}\|x\|^2$ is convex.

Proposition 11. If $f_t(x)$ is m_t -convex, $\int_{t \sim \mu} f_t(x)$ is $\int_{t \sim \mu} m_t$ -convex.

Proposition 12. If $f_t(x)$ is m_t -convex, $\sup_{t \in T} f_t(x)$ is $\inf_{t \in T} m_t$ -convex.

Definition 68. f is **strictly-convex** iff for all $x, y \in \text{dom } f$ such that $x \neq y$ and all $t \in (0, 1)$:

$$f(tx + (1 - t)y) < tf(x) + (1 - t)f(y) \quad (2.31)$$

Definition 69. f is **quasi-convex** iff for all $x, y \in \text{dom } f$ and all $t \in [0, 1]$:

$$f(tx + (1 - t)y) \leq \max\{f(x), f(y)\} \quad (2.32)$$

Definition 70. f is **pseudo-convex** iff for all $x, y \in \text{dom } f$:

$$(D^+ f)(x)(y - x) \geq 0 \implies f(y) \geq f(x) \quad (2.33)$$

Definition 71. f is **monotone** iff for all $x, y \in \text{dom } f$:

$$\langle f(x) - f(y), x - y \rangle \geq 0 \quad (2.34)$$

Definition 72. f is **m -monotone** iff $x \mapsto f(x) - mx$ is monotone. Equivalently, for all $x, y \in \text{dom } f$:

$$\langle f(x) - f(y), x - y \rangle \geq m\|x - y\|^2 \quad (2.35)$$

Definition 73. f is **β -cocoercive** iff for all $x, y \in \text{dom } f$:

$$\langle f(x) - f(y), x - y \rangle \geq \beta\|f(x) - f(y)\|^2 \quad (2.36)$$

623 **Definition 74.** f is (α, C) -Hölder continuous iff for all $x, y \in \text{dom } f$:

$$d(f(x), f(y)) \leq Cd(x, y)^\alpha \quad (2.37)$$

624 **Definition 75.** f is C -Lipschitz continuous iff it is $(1, C)$ -Hölder continuous.

625 **Definition 76.** Let (X, Σ, μ) be a measure space. Let $0 < p < \infty$. The L^p quasi-norm of a function
626 f is

$$\|f\|_p = \left(\int_\mu |f|^p \right)^{\frac{1}{p}} \quad (2.38)$$

627 It is a norm for $p \geq 1$.

628 **Definition 77.** Likewise, the L^∞ norm is the essential supremum of $|f|$:

$$\|f\|_\infty = \text{ess sup } |f| \quad (2.39)$$

629 Functions with finite L^∞ are also called **essentially bounded**.

630 **Definition 78.** $L^p(A; B)$ is the set of (equivalence classes of) functions from A to B with finite L^p
631 norm.

632 2.8 Pseudo-gradients

633 **Definition 79.** A **pseudo-gradient** is an estimator for the gradient of some smoothed approxima-
634 tion of a function. It uses only function evaluations, as opposed to exact analytical methods like
635 forward or reverse automatic differentiation.

636 A pseudo-gradient is typically used when the true gradient is unavailable, difficult to calculate,
637 or when the function is non-differentiable. Consider the problem of maximizing $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with
638 access to its values but not its derivatives. This setting is called **zeroth-order optimization**.
639 One approach to this problem computes estimates of the gradient $g(x) \approx \nabla f(x)$ and applies
640 gradient-based optimization. The gradient can be estimated via **finite differences** as

$$g(x)_i = \frac{1}{\sigma} (f(x + \sigma e_i) - f(x)) \quad (2.40)$$

641 for all $i \in [d]$, where e_i is the i -th standard basis vector and σ is a small scalar. However, the number
642 of queries needed scales linearly with the number of dimensions d .

643 To circumvent this scaling issue, an alternative approach evaluates the function at **randomly-**
644 **sampled points** and estimates the gradient as a sum of directional derivatives along **random**
645 **directions** (Duchi, Jordan, et al., 2015; Nesterov and Spokoiny, 2017; Shamir, 2017; Berahas et al.,
646 2022). These methods compute an unbiased estimator of the gradient of a *smoothed* version of f ,
647 which is induced by stochastically perturbing the input under a distribution μ_1 and taking the
648 expectation (Duchi, Bartlett, and Wainwright, 2012). Specifically, these take the form of

$$\nabla_x \mathbb{E}_{z \sim \mu_1} f(x + \sigma z) = \frac{1}{\sigma} \mathbb{E}_{z \sim \mu_2} f(x + \sigma z) z \quad (2.41)$$

649 where μ_1 and μ_2 are two (possibly different) distributions. For example, one can use the following
650 **distributions**:

- 651 • **Gaussian smoothing:** $\mu_1 = \mu_2 = \text{Normal}(0, I_d)$.
- 652 • **Ball smoothing:** $\mu_1 = \text{Uniform}(X)$ and $\mu_2 = \text{Uniform}(\partial X)$ where X is the d -dimensional unit
- 653 ball of radius \sqrt{d} .

654 In practice, we often estimate $\nabla_x \mathbb{E}_{z \sim \mu_1} f(x + \sigma z)$ by averaging *multiple* samples together:

$$\frac{1}{\sigma N} \sum_{i=1}^N a_i z_i \tag{2.42}$$

655 where $z_i \sim \mu_2$ are independent samples and a_i is evaluated using one of the following **stencils**:

- 656 • **Single-point:** $f(x + \sigma z_i)$
- 657 • **Forward-difference:** $f(x + \sigma z_i) - f(x)$
- 658 • **Central-difference:** $\frac{1}{2}(f(x + \sigma z_i) - f(x - \sigma z_i))$

659 Because the single-point stencil suffers from a large variance that diverges to infinity as σ approaches

660 0, the latter two are typically preferred in practice (Berahas et al., 2022).

661 These mathematical foundations underpin several classes of algorithms. Black-box zeroth-order

662 optimization uses only function evaluations to optimize a black-box function with respect to a set

663 of inputs. In particular, it does not require gradients. A prominent class of these algorithms is

664 **evolution strategies** (ES) (Rechenberg, 1973; Schwefel, 1977; Rechenberg, 1978; Bäck, 1996; Bäck,

665 Fogel, and Michalewicz, 1997; Eiben and Smith, 2003). Historically, these methods maintained and

666 evolved a *population* of parameter vectors.

667 **Natural evolution strategies** (NES) (Wierstra et al., 2014; Yi et al., 2009) represented the

668 population as a **distribution over parameters** and maximized its average objective value using

669 the score function estimator. For many parameter distributions, such as Gaussian smoothing, this is

670 equivalent to evaluating the function at randomly-sampled points and estimating the gradient as a

671 sum of estimates of directional derivatives along random directions (Fu, 2015; Duchi, Jordan, et al.,

672 2015; Nesterov and Spokoiny, 2017; Shamir, 2017; Berahas et al., 2022).

673 Salimans et al. (2017) applied Gaussian smoothing to single-agent reinforcement learning and

674 obtained competitive results on standard benchmarks. Lenc et al. (2019) showed that ES is a viable

675 method for learning non-differentiable parameters of large supervised models.

676 A closely related technique for smoothed gradient estimation is **simultaneous-perturbation**

677 **stochastic approximation** (SPSA). Spall (1992) introduced this method, which perturbs each

678 coordinate with Rademacher variates $\mu_1 = \mu_2 = \text{Uniform}(\{-1, 1\}^d)$ and uses the central-difference

679 stencil. A Taylor expansion of f shows that this is a good estimate of the true gradient when

680 σ is small. Later, Spall (1997) introduced a one-measurement variant of SPSA that utilized the

681 single-point stencil.

682 Chapter 3

683 Game theory definitions

684 In this section, we define some fundamental concepts in game theory. These include various types of
685 games, as well as solution concepts and performance metrics. We also give a mathematically rigorous
686 formulation of the setting we are tackling.

687 3.1 Normal-form game

688 von Neumann (1928) introduced the concept of a normal-form game. The motivation was to eliminate
689 the chaotic complexities of time, sequence, and hidden information in strategic interactions. Von
690 Neumann later expanded on this concept, and game theory in general, alongside economist Oskar
691 Morgenstern in their seminal book, *Theory of Games and Economic Behavior* (von Neumann and
692 Morgenstern, 1947).

693 **Definition 80.** A **normal-form game** (NFG) is a tuple (I, S, u) whose elements are as follows.

- 694 • I is a set of **players**.
- 695 • S_i is a set of **strategies** for each player i . A **strategy profile** is a map from each player to an
696 element of its strategy set. That is, it is an element of ΠS .
- 697 • $u \in \Pi S \times I \rightarrow \mathbb{R}$ is a **utility function**. It maps a strategy profile and player to a utility.

698 **Definition 81.** A **best response** (BR) is a strategy that maximizes a player's utility given the
699 other players' strategies. The set of best responses to strategy profile s by player i is denoted by

$$\text{BR}(s)_i = \operatorname{argmax}_{x \in S_i} u(s[i \mapsto x])_i \quad (3.1)$$

700 **Definition 82.** A **Nash equilibrium** (NE) is a strategy profile for which each player's strategy is
701 a best response to the other players' strategies. Thus

$$s \in \text{NE} \iff \forall i \in I, s_i \in \text{BR}(s)_i \quad (3.2)$$

702 **Definition 83.** A player's **BR value** is the supremum utility it could attain by unilaterally changing
703 its strategy. It is denoted by

$$\text{BRV}(s)_i = \sup_{x \in S_i} u(s[i \mapsto x])_i \quad (3.3)$$

704 **Definition 84.** A player’s **regret** is the supremum utility it could gain by unilaterally changing its
705 strategy. It is denoted by

$$\text{Reg}(s)_i = \text{BRV}(s)_i - u(s)_i \quad (3.4)$$

706 **Definition 85.** An ε -**Nash equilibrium** (ε -NE) is a strategy profile for which each player’s regret
707 is at most ε . Thus

$$s \in \varepsilon\text{-NE} \iff \sup_{i \in I} \text{Reg}(s)_i \leq \varepsilon \quad (3.5)$$

708 The special case $\varepsilon = 0$ coincides with NE.

709 In the following, let μ be a measure on I .

710 **Definition 86.** The **social welfare** of a strategy profile is the aggregate utility across players. It is
711 denoted by

$$\text{Wel}(s) = \int_{i \sim \mu} u(s)_i \quad (3.6)$$

712 **Definition 87.** The **exploitability** of a strategy profile is the aggregate regret across players. It is
713 denoted by

$$\text{Expl}(s) = \int_{i \sim \mu} \text{Reg}(s)_i \quad (3.7)$$

714 Exploitability is a standard metric of “closeness” to NE in the literature (Lanctot, Zambaldi,
715 et al., 2017; Lockhart et al., 2019; Walton and Lisy, 2021; Timbers et al., 2022). It is sometimes
716 called “NashConv”.

717 **Definition 88.** The **Nikaido–Isoda function** (NI) is

$$\text{NI}(x, y) = \int_{i \sim \mu} u(x[i \mapsto y]_i) - u(x)_i \quad (3.8)$$

718 The Nikaido–Isoda function was introduced by Nikaido and Isoda (1955), and studied by Flåm
719 and Antipin (1996), Flåm and Ruszczyński (2008), and Hou, Wen, and Chang (2018).

720 Note that

$$\text{Expl}(x) = \sup_{y \in \text{IS}} \text{NI}(x, y) \quad (3.9)$$

721 Therefore, finding an NE (if one exists) is equivalent to solving the min-max problem $\inf_x \sup_y \text{NI}(x, y)$.
722 Some prior work has used this function to search for NE (Uryasev and Rubinstein, 1994; Berridge
723 and Krawczyk, 1997; Krawczyk and Uryasev, 2000; Krawczyk, 2005; Flåm and Ruszczyński, 2008;
724 Gürkan and Pang, 2009; Heusinger and Kanzow, 2009a; Heusinger and Kanzow, 2009b; Qu and
725 Zhao, 2013; Hou, Wen, and Chang, 2018; Raghunathan, Cherian, and Jha, 2019; Tsaknakis and
726 Hong, 2021).

727 **Definition 89.** An **exact potential** is a function $\phi : \text{IS} \rightarrow \mathbb{R}$ such that

$$u(s[i \mapsto x]_i) - u(s)_i = \phi(s[i \mapsto x]) - \phi(s) \quad (3.10)$$

728 for all $s \in \text{IS}, i \in I, x \in S_i$. That is, a deviation increases utility by the amount that it increases
729 the potential.

730 **Definition 90.** An **ordinal potential** is a function $\phi : \Pi S \rightarrow \mathbb{R}$ such that

$$u(s[i \mapsto x])_i > u(s)_i \iff \phi(s[i \mapsto x]) > \phi(s) \quad (3.11)$$

731 for all $s \in \Pi S, i \in I, x \in S_i$. That is, a deviation increases utility iff it increases the potential.

732 **Proposition 13.** A **global maximum of a potential** is an NE.

733 *Proof.* Suppose $s \in \operatorname{argmax} \phi$. For every $i \in I$ and $x \in S_i$, $\phi(s[i \mapsto x]) \leq \phi(s)$, and thus $u(s[i \mapsto$
734 $x])_i \leq u(s)_i$. Thus, for every $i \in I$, $s_i \in \operatorname{BR}(s)_i$. Thus $s \in \operatorname{NE}$. \square

735 **Proposition 14.** If the utility function is **twice-differentiable**, an exact potential function exists
736 iff

$$\frac{\partial^2 u(s)_i}{\partial s_i \partial s_j} = \frac{\partial^2 u(s)_j}{\partial s_i \partial s_j} \quad (3.12)$$

737 for all $i, j \in I$. That is, the marginal effect of a player's strategy on another player's marginal utility
738 equals the marginal effect of the latter's strategy on the former. This theorem is due to Monderer
739 and Shapley (1996).

740 **Definition 91.** Let V be a topological vector space. Suppose $\forall i \in I, S_i \subseteq V$. An **aggregative**
741 **game** (Selten, 1970; Cornes and Hartley, 2007) is a game such that there exists a measure μ on I
742 and functions $f_i : S_i \times V \rightarrow \mathbb{R}$ for every $i \in I$ such that

$$\bar{s} = \int_{j \sim \mu} s_j \quad (3.13)$$

$$u(s)_i = f_i(s_i, \bar{s}) \quad (3.14)$$

743 for every $i \in I$. Here, \bar{s} is the **aggregate strategy** across players. In other words, each player's
744 utility depends on its own strategy and on an aggregate of all players' strategies.

745 3.2 Mixed-strategy game

746 **Definition 92.** A **mixed-strategy game** is a tuple (I, A, R) whose elements are as follows.

- 747 • I is a set of **players**.
- 748 • A_i is a set of **actions** for each player i .
- 749 • $R : \Pi A \rightarrow \mathbb{R}$ is a **reward function**.

750 **Definition 93.** A **mixed strategy** for player i is a probability measure $\sigma_i \in \Delta A_i$.

751 **Definition 94.** The **joint mixed strategy** is the product measure $\sigma \in \Delta \Pi A$ defined by

$$\sigma(a) = \bigotimes_{i \in I} \sigma_i(a_i) \quad (3.15)$$

752 **Definition 95.** The **utility** is the expected reward under sampling of an action profile:

$$u(\sigma) = \mathbb{E}_{a \sim \sigma} R(a) \quad (3.16)$$

753 **Definition 96.** The **induced NFG** is $(I, \prod_{i \in I} \Delta A_i, u)$. That is, each player chooses a mixed
 754 strategy, and the utility is the expected utility under the resulting joint mixed strategy.

755 **Definition 97.** A **mixed strategy Nash equilibrium** (MSNE) of the NFG (I, A, R) is an NE of
 756 the aforementioned induced NFG.

757 **Definition 98.** A **pure strategy Nash equilibrium** (PSNE) is an MSNE where each factor
 758 distribution is a Dirac delta function.

759 **Proposition 15.** A player's BR value can be computed by maximizing over actions:

$$\text{BRV}(\sigma)_i = \sup_{x \in \Delta A_i} u(\sigma[i \mapsto x])_i = \sup_{x \in A_i} u(\sigma[i \mapsto \delta_x])_i \quad (3.17)$$

760 **Definition 99.** A **correlated equilibrium** (CE) (Aumann, 1974) is a distribution $\sigma \in \Delta \Pi A$ over
 761 action profiles such that

$$\mathbb{E}_{a \sim \sigma} u(a)_i \geq \mathbb{E}_{a \sim \sigma} u(a[i \mapsto f_i(a_i)])_i \quad (3.18)$$

762 for every player $i \in I$ and deviation function $f_i : A_i \rightarrow A_i$. Unlike NE, which has to be a product
 763 distribution, this distribution can be correlated.

764 **Definition 100.** A **coarse correlated equilibrium** (CCE) (Moulin and Vial, 1978) is a distribution
 765 $\sigma \in \Delta \Pi A$ over action profiles such that

$$\mathbb{E}_{a \sim \sigma} u(a)_i \geq \mathbb{E}_{a \sim \sigma} u(a[i \mapsto a'_i])_i \quad (3.19)$$

766 for every player $i \in I$ and deviation action $a'_i \in A_i$.

767 **Proposition 16.** The previous solution concepts are related as follows.

$$\text{PSNE} \subseteq \text{MSNE} \subseteq \text{CE} \subseteq \text{CCE} \quad (3.20)$$

768 3.3 Bayesian game

769 Harsanyi (1967) introduced the concept of a Bayesian game. The motivation was to formalize the
 770 concept of **incomplete information**, where players are uncertain about the rules, payoffs, or
 771 others' available moves. **Harsanyi's transformation** converts an incomplete information game
 772 into a complete but imperfect information game. It does this as follows. There is a move by Nature
 773 at the start of the game. Nature randomly assigns each player a "type" or observation, representing
 774 their private information (e.g., cost or preference). Players know their own type but only have a
 775 probability distribution (beliefs) over others' types. By framing the uncertainty as a move by Nature,
 776 this allowed existing techniques and tools for NE to be used.

777 **Definition 101.** A **Bayesian game** is a tuple (I, S, O, A, Z, R, ρ) whose elements are as follows.
 778 We assume all sets are equipped with σ -algebras (as described in Section 2.5.1) and all functions are
 779 measurable.

- 780 • I is a set of **players**.

- 781 • S is a set of **states**.
- 782 • O_i is a set of **observations** for each player i .
- 783 • A_i is a set of **actions** for each player i .
- 784 • $Z : S \rightarrow \Delta \Pi O$ is an **observation function**.
- 785 • $R : S \times \Pi A \rightarrow \mathbb{R}^I$ is a **reward function**.
- 786 • $\rho \in \Delta S$ is a **state distribution**.

787 **Definition 102.** A **policy** for player i is a function $\pi_i : O_i \rightarrow \Delta A_i$. It maps each observation to a
788 distribution of actions.

789 **Definition 103.** The **joint policy** is the function $\pi : \Pi O \rightarrow \Delta \Pi A$ defined by

$$\pi(a | o) = \bigotimes_{i \in I} \pi_i(a_i | o_i) \quad (3.21)$$

790 **Definition 104.** The **utility** is the expected reward under the following process:

$$u(\pi) = \mathbb{E}_{s \sim \rho} \mathbb{E}_{o \sim Z(s)} \mathbb{E}_{a \sim \pi(o)} \mathbb{E}_{r \sim R(s,a)} r \quad (3.22)$$

791 That is, a state is sampled, a joint observation is sampled given the state, a joint action is sampled
792 given the observation, and a joint reward is sampled given the state and joint action.

793 **Proposition 17.** A player's BR value can be computed as follows:

$$\text{BRV}(\pi)_i = \sup_{x: O_i \rightarrow \Delta A_i} u(\pi[i \mapsto x])_i \quad (3.23)$$

$$= \sup_{x: O_i \rightarrow \Delta A_i} \mathbb{E}_{s \sim \rho} \mathbb{E}_{o \sim Z(s)} \mathbb{E}_{a \sim \pi[i \mapsto x](o)} \mathbb{E}_{r \sim R(s,a)} r_i \quad (3.24)$$

$$= \mathbb{E}_{o_i} \sup_{x \in \Delta A_i} \mathbb{E}_{a_i \sim x} \mathbb{E}_{s | o_i} \mathbb{E}_{a_{-i} | s} \mathbb{E}_{r \sim R(s,a)} r_i \quad (3.25)$$

$$= \mathbb{E}_{o_i} \sup_{a_i \in A_i} \mathbb{E}_{s | o_i} \mathbb{E}_{a_{-i} | s} \mathbb{E}_{r \sim R(s,a)} r_i \quad (3.26)$$

794 This means that, instead of maximizing over the policy space $O_i \rightarrow \Delta A_i$, we can sample observations,
795 maximize over the action space A_i , and sample states conditioned on the observation. This technique
796 was used by Bichler, Fichtl, Heidekrüger, et al. (2021, Supplementary Section 3) to estimate
797 exploitabilities.

798 3.4 Mean-field game

799 **Definition 105.** A **mean-field game** (MFG) is a tuple $(X, A, T, R, \rho, \gamma)$ whose elements are as
800 follows. We assume that all sets are equipped with σ -algebras (as described in Section 2.5.1) and all
801 functions are measurable.

- 802 • X is a set of **states** that an individual player can occupy.
- 803 • A is a set of **actions** that an individual player can perform.

- 804 • $T : X \times M(X) \times A \rightarrow \Delta X$ is a **transition function**.
- 805 • $R : X \times M(X) \times A \rightarrow \mathbb{R}$ is a **reward function**.
- 806 • $\rho \in M(X)$ is an **initial state distribution**.
- 807 • $\gamma \in \mathbb{R}$ is a **discount factor**.

808 Here, $M(X)$ is the set of measures on X . A **population distribution** is an element of $M(X)$,
 809 representing a distribution of players over the set of states X , i.e., an occupancy measure for the
 810 states.

811 A **Markov policy** is a sequence $\{\pi_t\}_{t \in \mathbb{N}}$ of stochastic kernels on A given X .

812 A **population flow** is a sequence $\{m_t\}_{t \in \mathbb{N}}$ of measures on X .

813 If a player executes a Markov policy π under an exogenous population flow m , its individual
 814 **trajectory** is generated as follows:

$$x_0 \sim \rho \tag{3.27}$$

$$a_t \sim \pi_t(x_t) \tag{3.28}$$

$$r_t = R(x_t, m_t, a_t) \tag{3.29}$$

$$x_{t+1} \sim T(x_t, m_t, a_t) \tag{3.30}$$

815 The player's **expected return** (cumulative reward) is

$$J(\pi, m) = \mathbb{E} \sum_{t=0}^{\infty} \gamma^t r_t \tag{3.31}$$

816 Equivalently,

$$q_t(x, a) = R(x, m_t, a) + \gamma \mathbb{E}_{x' \sim T(x, m_t, a)} v_{t+1}(x') \tag{3.32}$$

$$v_t(x) = \mathbb{E}_{a \sim \pi_t(x)} q_t(x, a) \tag{3.33}$$

$$J(\pi, m) = \mathbb{E}_{x \sim \rho} v_0(x) \tag{3.34}$$

817 If a **representative player** follows a Markov policy π , the **induced population flow** $m = \Lambda(\pi)$
 818 is

$$m_0 = \rho \tag{3.35}$$

$$m_{t+1}(B) = \int_{x_t \sim m_t} \mathbb{E}_{a_t \sim \pi_t(x_t)} T(x_t, m_t, a_t)(B) \tag{3.36}$$

819 A pair (π, m) consisting of a Markov policy π and population flow m is a **mean-field equilibrium**
 820 iff the following two conditions hold:

$$\pi \in \operatorname{argmax}_{\pi} J(\pi, m) \tag{3.37}$$

rationality

$$m = \Lambda(\pi) \tag{3.38}$$

self-consistency

3.5 Partially observable stochastic game

Definition 106. A **partially observable stochastic game** (POSG) is a tuple $(I, S, O, A, Z, T, V, R, \Gamma, \rho)$ whose elements are as follows. We assume that all sets are equipped with σ -algebras (as described in Section 2.5.1) and all functions are measurable.

- I is a set of **players**.
- S is a set of **states**.
- O_i is a set of **observations** for each player i .
- A_i is a set of **actions** for each player i .
- $Z : S \rightarrow \Delta \Pi O$ is an **observation function**.
- $T : S \times \Pi A \rightarrow \Delta S$ is a **transition function**.
- V is a Banach space of **value vectors**, which are functions in $I \rightarrow \mathbb{R}$ that are essentially bounded (c.f. Section 2.7) and measurable (c.f. Section 2.5.1).
- $R : S \times \Pi A \times S \rightarrow \Delta V$ is a **reward function**.
- $\Gamma : S \times \Pi A \times S \rightarrow \Delta \mathcal{L}(V)$ is a **discount function**.
- $\rho \in \Delta S$ is an **initial state distribution**.

Notice that a discount factor is a linear operator on the space of value vectors.¹

Definition 107. A **policy** for player i is a function $\pi_i : O_i \rightarrow \Delta A_i$. It maps each observation to a distribution of actions. For simplicity of exposition, the policies are assumed to be *reactive*, that is, depend only on the current observation. This models imperfect recall games with no memory. However, this can be extended to policies *with* memory.

¹Many linear operators can be expressed as **kernels**, that is, as integrals. However, in infinite spaces, there are more exotic operators that cannot be expressed in this way. An example is the limit operator. Suppose the set of players is \mathbb{N} . Consider the linear operator γ such that Player 0's value depends on the long-term trend of the entire population, regardless of any finite group of individuals:

$$(\gamma v)(0) = \lim_{n \rightarrow \infty} v(n) \tag{3.39}$$

(Technically, since the limit does not always exist, we extend this to a **Banach limit**, which is a non-negative, shift-invariant, linear functional that extends the usual limit and assigns a limit-like value to all bounded sequences.) This is a linear operator, but it cannot be represented as a kernel. Conceptually, the “mass” of this “measure” is not located on any specific integer, nor spread out as a density. Instead, its mass is concentrated “at infinity”—more precisely, in the so-called **Stone-Ćech compactification** of \mathbb{N} . Informally, the Stone-Ćech compactification of a topological space captures the “ways of going to infinity” in that space and adjoins them to it.

Kernel operators cover most cases in physics and economics (local interactions, weighted averages, and mean fields). Exotic operators (like the Banach limit) represent “ghost” interactions with the infinite tail of the population that cannot be localized. More precisely, the Yosida-Hewitt decomposition (Yosida and Hewitt, 1952) shows that a finitely-additive measure can be uniquely decomposed into the sum of a countably additive measure (the kernel component, in which the mass is distributed across points or “small” sets) and a “purely” finitely additive measure (the exotic component, in which the mass exists at the “boundary” or “infinity”).

Exotic operators can model esoteric phenomena like “rational bubbles” in infinite economies where value comes from “infinity” rather than any fundamental agent (Bewley, 1972; Gilles and LeRoy, 1992). Under the **Yosida-Hewitt decomposition**, the price of an asset over an infinite horizon decomposes into the sum of the **fundamental value** and the **rational bubble**. The fundamental value is the standard discounted sum of all future dividends paid to specific, localized agents in time. The rational bubble is a component of the price that exists purely because agents believe they can sell the asset to someone else in the future, *ad infinitum*. It has been argued that **fiat money** is an example of a rational bubble (Gilles and LeRoy, 1992; Santos and Woodford, 1997; Huang and Werner, 2000).

841 **Definition 108.** The **joint policy** is the function $\pi : \Pi O \rightarrow \Delta \Pi A$ defined by

$$\pi(a \mid o) = \bigotimes_{i \in I} \pi_i(a_i \mid o_i) \quad (3.40)$$

842 3.5.1 Episode formulation

843 **Definition 109.** An **episode** is a tuple (s, o, a, r, γ) where $s_0 \sim \rho$ and, for each timestep $t \in \mathbb{N}$,

$$o_t \sim Z(s_t) \quad (3.41)$$

$$a_t \sim \pi(o_t) \quad (3.42)$$

$$s_{t+1} \sim T(s_t, a_t) \quad (3.43)$$

$$r_t \sim R(s_t, a_t, s_{t+1}) \quad (3.44)$$

$$\gamma_t \sim \Gamma(s_t, a_t, s_{t+1}) \quad (3.45)$$

844 **Definition 110.** The **return** for timestep t is

$$g_t = r_t + \gamma_t g_{t+1} \quad (3.46)$$

845 **Definition 111.** The **cumulative discount** from timestep i to timestep $j \geq i$ is

$$\rho_{ij} = \prod_{k=i}^{j-1} \gamma_k \quad (3.47)$$

846 Thus the return for timestep t is

$$g_i = \sum_{j=i}^{\infty} \rho_{ij} r_j \quad (3.48)$$

847 **Definition 112.** The **utility** is

$$u(\pi) = \mathbb{E} g_0 \quad (3.49)$$

848 **Definition 113.** The **induced NFG** is $(I, \prod_{i \in I} (O_i \rightarrow \Delta A_i), u)$. That is, each player chooses a
849 policy, and the utility is the expected utility under the resulting joint policy.

850 3.5.2 Value formulation

851 An equivalent way to formulate the expected utility is as follows. Given a joint policy π , let

- 852 • $q(s, a)$ be the value of state s and action a .
- 853 • $v(s)$ be the value of state s .
- 854 • $u(\pi)$ be the utility.

855 Then

$$q(s, a) = \mathbb{E}_{s' \sim T(s, a)} \mathbb{E}_{r \sim R(s, a, s')} \mathbb{E}_{\gamma \sim \Gamma(s, a, s')} r + \gamma v(s') \quad (3.50)$$

$$v(s) = \mathbb{E}_{o \sim Z(s)} \mathbb{E}_{a \sim \pi(o)} q(s, a) \quad (3.51)$$

$$u(\pi) = \mathbb{E}_{s \sim \rho} v(s) \quad (3.52)$$

856 where q and v are defined recursively in terms of each other.

857 Chapter 4

858 Prior work

859 In this section, we describe prior work in each of the main areas we are tackling.

860 4.1 Prior work on infinite-action games

861 Ganzfried and Sandholm (2010) presented a procedure for solving large imperfect-information games
862 by solving an infinite approximation of the original game and mapping the equilibrium to a strategy
863 profile in the finite game. Perhaps counterintuitively, this infinite approximation could often be
864 solved much more easily than the original finite game. To find an equilibrium in continuous games,
865 the algorithm exploited a qualitative model of the equilibrium structure as an additional input.

866 4.1.1 Action space discretization

867 The typical approach to computing an equilibrium of a game with continuous action spaces involves
868 discretizing the action space. Building on this approach, Bichler, Fichtl, and Oberlechner (2023)
869 introduced an algorithmic framework for Bayesian games with continuous type and action spaces.
870 Specifically, they discretized the type and action spaces and implemented gradient dynamics in
871 the discretized version of the game without using neural networks. They computed distributional
872 strategies (Milgrom and Weber, 1985) (a form of mixed strategies for Bayesian games) via online
873 convex optimization, namely **simultaneous online dual averaging** (SODA). Consequently, they
874 demonstrated that the equilibrium of their discretized game approximated an equilibrium in the
875 continuous game.

876 Historically, a common approach to handling continuous action spaces has been to **discretize**
877 them. Although this method can work well for smaller games, it inherently entails a loss in solution
878 quality. For example, Kroer and Sandholm (2015) provided bounds on this quality degradation for
879 extensive-form games. Furthermore, this approach fails to scale to high-dimensional observation and
880 action spaces, because multidimensional environments cause a combinatorial explosion of points
881 that scales exponentially with the number of dimensions. This scalability bottleneck is particularly
882 intractable in highly complex spaces, such as generative adversarial networks (GANs) (Goodfellow
883 et al., 2014), where a strategy is an entire probability distribution over images generated by a neural
884 network. Additionally, explicit gradient information about the game is often unavailable in many
885 practical applications. These limitations highlight a strong need for alternative approaches.

886 4.1.2 Double oracle

887 McMahan, Gordon, and Blum (2003) introduced the **double oracle** (DO) algorithm for normal-form
888 games and proved its convergence. This algorithm maintains finite strategy sets for each player and
889 iteratively expands them with best responses to an equilibrium of the induced finite sub-game. Adam
890 et al. (2021) extended DO to two-player zero-sum continuous games, and Kroupa and Votroubek
891 (2023) generalized it to n -player continuous games. This approach converged fast when the strategy
892 space dimension was small and the generated subgames were not excessively large. For example,
893 in two-player zero-sum cases, best responses were computed by selecting the optimal point of a
894 uniform discretization for one-dimensional problems or by using a mixed-integer linear programming
895 reformulation for Colonel Blotto games. However, this approach did not scale well to high-dimensional
896 games with general payoffs where best-response computation was difficult. Moreover, estimating the
897 finite subgame for stochastic games proved challenging and required many samples. Furthermore,
898 this method failed to learn observation-dependent strategies that generalized across observations.

899 Li and Wellman (2021) extended the double oracle approach to n -player general-sum continuous
900 Bayesian games by representing agents as neural networks optimized via NES (Wierstra et al.,
901 2014). To approximate a pure-strategy equilibrium, they formulated the problem as a bi-level
902 optimization and employed NES for both inner-loop best-response optimization and outer-loop
903 regret minimization.

904 Double oracle algorithms generally maintained a set of static strategies that expanded on each
905 iteration with approximate best responses to the opponents’ meta-strategies. Several prominent
906 algorithms built upon this foundational concept, including PSRO and its variants, which are described
907 below.

908 Lanctot, Zambaldi, et al. (2017) introduced **policy-space response oracles** (PSRO), which
909 generalized DO by defining the metagame’s choices as continuous policies rather than discrete
910 actions. PSRO successfully generalized fictitious self-play and allowed any meta-solver to compute
911 new meta-strategies.

912 McAleer, Lanier, Fox, et al. (2020) introduced **pipeline PSRO** (P2SRO), a scalable method for
913 finding approximate Nash equilibria (NE) in large zero-sum imperfect-information games. P2SRO
914 parallelized PSRO with convergence guarantees by maintaining a hierarchical pipeline of reinforcement
915 learning workers that trained against policies generated by lower hierarchy levels.

916 McAleer, Lanier, Wang, Baldi, and Fox (2021) introduced **extensive-form DO** (XDO) for
917 two-player zero-sum games, which guaranteed linear convergence to an approximate NE relative
918 to the number of infostates. Unlike PSRO, which mixed best responses at the root of the game,
919 XDO mixed best responses at every infostate. The authors also proposed **neural XDO** (NXDO),
920 which learned best responses through deep reinforcement learning and computed the metagame
921 equilibrium using methods like NFSP (Heinrich and Silver, 2016) or DREAM (Steinberger, Lerer, and
922 Brown, 2020). NXDO iteratively added reinforcement learning policies to a population while solving
923 an extensive-form restricted game, a process that proved more efficient than solving matrix-form
924 restricted games.

925 McAleer, Wang, et al. (2022) developed **anytime DO** (ADO), a tabular algorithm for two-player
926 zero-sum games that guaranteed convergence to an NE while monotonically decreasing exploitability
927 across iterations. They additionally introduced **anytime PSRO** (APSRO), an extension of ADO
928 that calculated best responses using reinforcement learning.

929 McAleer, Lanier, Wang, Baldi, Sandholm, et al. (2024) created **self-play PSRO** (SP-PSRO), a
930 method restricted to two-player zero-sum games that added approximately optimal stochastic policies
931 to the population in each iteration. Instead of exclusively adding deterministic best responses to

932 the opponent’s least exploitable population mixture, SP-PSRO learned and incorporated stochastic
933 policies to avoid the need to enumerate all deterministic policies prior to convergence. However,
934 SP-PSRO remained a normal-form algorithm that mixed at the root of the game tree, which McAleer,
935 Lanier, Wang, Baldi, and Fox (2021) previously showed could require an exponential number of
936 iterations to converge. The authors left the combination of SP-PSRO with extensive-form techniques
937 like XDO and NXDO for future work.

938 Muller, Omidshafiei, et al. (2020) extended the theoretical underpinnings of PSRO by evaluating
939 an alternative solution concept called α -Rank (Omidshafiei et al., 2019) instead of NE. They
940 established convergence guarantees for this approach and identified formal links between NE and
941 α -Rank.

942 Finally, Marris et al. (2021) proposed **joint policy-space response oracles** (JPSRO) for
943 training agents in n -player, general-sum extensive-form games. This algorithm shifted the solution
944 concept away from NE, utilizing CE and CCE instead (as described in Section 3.2).

945 4.1.3 Fictitious play

946 Brown (1951) introduced **fictitious play** (FP), a popular game-theoretic model of learning where
947 players repeatedly play a game and choose a best response to their opponents’ historical average
948 strategies at each iteration. The average strategy profile converges to a Nash equilibrium (NE)
949 in certain classes of games, including two-player zero-sum and potential games. Players update
950 their beliefs simultaneously or alternately (Berger, 2007). Leslie and Collins (2006) introduced
951 **generalized weakened FP** (GWFP), which generalizes FP by allowing approximate best responses
952 and perturbed average strategy updates.

953 Perkins and Leslie (2014) extended stochastic FP to the continuous action space framework.
954 They studied the limiting behavior of stochastic FP using the associated smooth best response
955 dynamics on the space of finite signed measures. Using this approach, they showed that stochastic
956 FP converges to an equilibrium in two-player zero-sum games.

957 Heinrich, Lanctot, and Silver (2015) introduced **full-width extensive-form FP** (XFP), which
958 extends FP to extensive-form multi-step games. They also introduced **fictitious self-play** (FSP), a
959 sample-based machine learning framework that implements GWFP in behavioral strategies. FSP
960 avoids cycles by computing an approximate best response against a uniform mixture of all previous
961 policies, which converges to an NE in two-player zero-sum games (Heinrich, Lanctot, and Silver,
962 2015). In FSP, players repeatedly play a game, store their experiences in memory, and mix between
963 their best responses and average strategies to act cautiously. At each iteration, players replay their
964 experience of play against opponents to compute an approximate best response, and similarly replay
965 their own behavioral experience to learn a model of their average strategy.

966 Heinrich and Silver (2016) introduced **neural fictitious self-play** (NFSP), which combines
967 FSP with neural network function approximation and deep reinforcement learning. An NFSP
968 agent consists of two neural networks. The first network is trained via reinforcement learning from
969 memorized experiences against fellow agents to learn an approximate best response to their historical
970 behavior. The second network is trained via supervised learning from memorized experiences of the
971 agent’s own behavior.

972 Kamra, Gupta, Wang, et al. (2019) introduced **DeepFP**, an approximate FP algorithm for two-
973 player games with continuous action spaces. They demonstrated stable convergence to equilibrium
974 on several classic games and a large forest security domain. DeepFP represents players’ approximate
975 best responses via generative neural networks, which act as highly expressive implicit density

976 approximators. Because these implicit density models cannot be trained directly in the absence of
977 gradients, the authors employed a game-model network to serve as a differentiable approximation of
978 the players’ utilities given their actions. This approach allows the networks to be trained end-to-end
979 in a model-based learning regime.

980 Ganzfried (2021) introduced **redundant fictitious play** (RFP), an algorithm for approximating
981 equilibria in continuous games, and applied it to a continuous Colonel Blotto game. Unlike our
982 approach, this algorithm requires a best-response oracle as a subroutine (e.g., a mixed-integer linear
983 program solver for the continuous Colonel Blotto game).

984 Vinyals et al. (2019) tackled StarCraft II, a real-time strategy game serving as a popular artificial
985 intelligence benchmark, by introducing a multiagent reinforcement learning algorithm called **league**
986 **training**. Noting that standard self-play algorithms learn rapidly but can chase cycles indefinitely,
987 they extended FSP to compute best responses against a non-uniform mixture of opponents. Their
988 league of potential opponents includes a diverse range of agents and their policies from both current
989 and previous iterations. At each iteration, agents play games against opponents sampled from an
990 agent-specific mixture policy, and their parameters are updated using an actor-critic reinforcement
991 learning procedure based on those game outcomes. The league consists of three distinct types of
992 agents that differ primarily in their mechanism for selecting the opponent mixture. First, main agents
993 utilize a **prioritized fictitious self-play** (PFSP) mechanism that adapts mixture probabilities
994 proportionally to each opponent’s win rate against the agent, providing more opportunities to
995 overcome problematic opponents. To recover the rapid learning of self-play, a main agent is selected
996 as an opponent with a fixed probability. Second, main exploiter agents play exclusively against
997 the current iteration of main agents to identify potential exploits and encourage the main agents
998 to address their weaknesses. Third, league exploiter agents use a PFSP mechanism similar to the
999 main agents, but are not targeted by main exploiters, aiming instead to find systemic weaknesses
1000 across the entire league. Both types of exploiter agents are periodically reinitialized to encourage
1001 diversity, allowing them to rapidly discover specialist strategies that are not necessarily robust
1002 against exploitation.

1003 4.1.4 Evolution strategies

1004 Bichler, Fichtl, Heidekrüger, et al. (2021) focused on symmetric auction models, assuming symmetric
1005 prior distributions and pure, symmetric equilibrium bidding strategies. To tackle these environments,
1006 they introduced a learning method called **neural pseudo-gradient ascent** (NPGA). NPGA
1007 represents strategies as neural networks and applies policy iteration based on gradient dynamics in
1008 self-play to provably learn local equilibria. The method follows the simultaneous gradient of the game
1009 and employs a smoothing technique to circumvent discontinuities in the *ex post* utility functions,
1010 which arise at the bid value where an arbitrarily small change makes the difference between winning
1011 and losing. They demonstrated that NPGA converges to a Bayesian Nash equilibrium across a wide
1012 variety of symmetric auction games.

1013 Building on this work, Bichler, Kohring, and Heidekrüger (2023) extended the approach to
1014 asymmetric auctions, a setting that requires training multiple neural networks. They analyzed a
1015 wide variety of asymmetric auction models and showed that their method closely approximates
1016 Bayesian Nash equilibria in all environments where analytical solutions are known. Furthermore,
1017 they evaluated new, larger environments lacking analytical solutions and verified that the discovered
1018 solutions still approximate the equilibrium closely.

1019 Li and Wellman (2021) tackled the problem of solving symmetric one-shot Bayesian games

1020 characterized by high-dimensional type and action spaces, multiple players, general-sum payoffs,
 1021 and no provided analytic structure. They represented agent strategies in parametric form as neural
 1022 networks and applied NES to optimize them. For pure equilibrium computation, they formulated
 1023 the problem as a bi-level optimization, using NES to implement both inner-loop best response
 1024 optimization and outer-loop regret minimization. For mixed equilibrium computation, they adopted
 1025 an incremental strategy generation framework where NES produces a finite sequence of approximate
 1026 best-response strategies. Equilibria are then calculated over this finite strategy set via a model-based
 1027 optimization process. Both methods use NES to search for strategies over the functional space of
 1028 policies, relying solely on black-box simulation access to noisy payoff samples.

1029 4.2 Prior work on infinite-player games

1030 This section reviews existing literature on solving infinite-player games. Sandholm (2001) studied
 1031 potential games with continuous player sets, such as random matching games with common payoffs
 1032 and congestion games. Parise and Ozdaglar (2019) and Parise and Ozdaglar (2023) presented a
 1033 framework for analyzing equilibria and designing interventions in large network games. These network
 1034 games are modeled using a graphon, where a continuum of players receives payoffs dependent on
 1035 their own action and a weighted average of others’ actions. Building on this, graphon mean-field
 1036 games (Caines and Huang, 2021; Cui and Koepl, 2022) consider heterogeneous players with a
 1037 continuous index space.

1038 Perrin et al. (2020) analyzed continuous-time fictitious play on finite-state MFGs with common
 1039 noise. They provided a convergence analysis and proved that the induced exploitability decreases
 1040 at a linear rate. Muller, Rowland, et al. (2022) introduced mean-field PSRO, an adaptation of
 1041 **policy-space response oracles** (PSRO) (Lanctot, Zambaldi, et al., 2017) that learns Nash, coarse-
 1042 correlated, and correlated equilibria in MFGs. Pérolat et al. (2022) applied **online mirror descent**
 1043 (OMD) to MFGs and proved that continuous-time OMD converges to a Nash equilibrium under
 1044 certain conditions. Wang and Wellman (2023) solved MFGs using a double oracle algorithm by
 1045 iteratively adding approximate best responses to the equilibrium of the empirical MFG formed by
 1046 previously considered strategies. Wu et al. (2024) proposed a deep reinforcement learning algorithm
 1047 that achieves population-dependent equilibria in MFGs without averaging or sampling from history.
 1048 Zhang, Chen, and Di (2025) introduced **SemiSGD**, a stochastic semi-gradient descent method
 1049 that treats a player’s policy and the population distribution as a unified parameter to allow fully
 1050 asynchronous, simultaneous updates. They also presented a **population-aware linear function**
 1051 **approximation** (PA-LFA) for continuous state-action MFGs.

1052 A parallel line of research applies reinforcement learning to MFGs under unknown dynamics. Guo
 1053 et al. (2019) proposed a fitted Q-iteration scheme that provably converges to mean-field equilibria in
 1054 discrete state-action MFGs by alternating between Q-value updates and a fixed-point update for
 1055 the population distribution. Xie et al. (2021) extended this program to the finite-horizon setting,
 1056 establishing sample-complexity bounds for provably efficient RL under contraction properties of
 1057 the mean-field operator. Laurière et al. (2022) surveyed the field, cataloging both model-based
 1058 and model-free approaches and highlighting open problems in non-stationary, non-contractive, and
 1059 continuous state-action MFGs.

1060 **Comparison of MFGs to our setting:** Classical MFG methods also address infinite-player
 1061 games, but rely on the fundamental assumption that players are indistinguishable and individually
 1062 negligible. Specifically, an individual’s payoff depends only on its own action and the anonymous

1063 **mean-field density** of the population.¹ In contrast, our approach does not require these symmetry
 1064 or mean-field assumptions. Each player can have an arbitrary strategy space and a unique utility
 1065 function that is not constrained to depend solely on aggregate behavior. Graphon MFGs (e.g.,
 1066 Caines and Huang (2021) and Cui and Koepl (2022) relax indistinguishability by indexing players
 1067 on a continuum, but still restrict each player’s payoff to depend on others’ actions only through a
 1068 fixed graphon-weighted aggregate. Our framework admits arbitrary measurable utility functionals,
 1069 including higher-order interactions and discontinuities, without a prescribed aggregation structure.

1070 **Agent-based modeling** (ABM) is a bottom-up computational method that simulates au-
 1071 tonomous, interacting agents, each with its own state and behavioral rules, within an environment
 1072 to observe how system-level patterns emerge from local interactions (Epstein and Axtell, 1996;
 1073 Bonabeau, 2002; Grimm et al., 2005; Macal and North, 2010). Our proposed framework acts as
 1074 a differentiable analogue to ABM, leveraging the generalization capabilities of neural networks to
 1075 accelerate learning.

1076 4.3 Prior work on learning dynamics

1077 In this section, we review prior work on gradient-based learning dynamics for games. They can
 1078 be characterized by the ordinary differential equations (ODEs) shown in Table 4.1. Here, \dot{x} is the
 1079 time derivative of the current strategy profile, and $v = \text{diag } \nabla u$ is the **simultaneous gradient**.
 1080 Each component $v(x)_i = \nabla_i u(x)_i$ is the gradient of a player’s utility with respect to their strategy.
 1081 Additionally, $J = \nabla v$ is the Jacobian of the vector field v , J^\top is its transpose, $J_a = \frac{1}{2}(J - J^\top)$ is its
 1082 antisymmetric part, and J_o is its off-diagonal part (replacing its diagonal with zeroes). Dots indicate
 1083 derivatives with respect to time, $\gamma > 0$ is a hyperparameter, and $v|_y$ denotes v evaluated at y rather
 1084 than x . These methods were extensively analyzed in prior work, including Balduzzi et al. (2018),
 1085 Letcher, Foerster, et al. (2019), Letcher, Balduzzi, et al. (2019), Mertikopoulos and Zhou (2019),
 1086 Grnarova, Levy, et al. (2019), Mazumdar, Sastry, and Jordan (2025), Hsieh, Mertikopoulos, and
 1087 Cevher (2021), and Willi et al. (2022).

1088 The simultaneous gradient yields a vector field on the space of strategy profiles. Because this
 1089 vector field may not be conservative (i.e., the gradient of a potential) as in gradient descent, standard
 1090 gradient-based optimization methods often struggle, and trajectories can cycle around fixed points
 1091 rather than converging to them. The actual optimization is done by discretizing each ODE in time.
 1092 For example, SG is discretized as $x_{i+1} = x_i + \eta v_i$ and OP is discretized as $x_{i+1} = x_i + \eta v_i + \gamma(v_i - v_{i-1})$,
 1093 where $\eta > 0$ is a stepsize and $v_j = v(x_j)$.

1094 **Simultaneous gradient ascent** (SG) maximizes each player’s utility independently, as if the
 1095 other players are fixed. In the two-player zero-sum case, it is also known as **gradient descent**
 1096 **ascent** (GDA). Goktas and Greenwald (2022b) studied min-max games with dependent strategy
 1097 sets, where the strategy of the first player constrains the behavior of the second. They introduced two
 1098 variants of GDA that assume access to a solution oracle for the optimal **Karush–Kuhn–Tucker**
 1099 (KKT) multipliers of the games’ constraints, proving a convergence guarantee.

1100 **Extragradient** (EG) (Korpelevich, 1976) takes a step in the direction of the simultaneous
 1101 gradient and uses the simultaneous gradient at that new point to take a step from the original point.
 1102 Golowich et al. (2020) proved a tight last-iterate convergence guarantee for EG.

1103 **Optimistic gradient** (OP) (Popov, 1980; Daskalakis et al., 2018; Hsieh, Iutzeler, et al., 2019)
 1104 uses past gradients to predict future gradients and updates according to the latter.

¹In anonymous games, identities do not matter, and only aggregate behavior affects outcomes.

Table 4.1: ODE corresponding to each method.

Method	\dot{x}
SG	v
EG	$v _{x+\gamma v}$
OP	$(I + \gamma \frac{d}{dt})v = v + \gamma \dot{v}$
CO	$(I - \gamma J^\top)v = v - \gamma \nabla \frac{1}{2} \ v\ ^2$
SGA	$(I - \gamma J_a^\top)v$
LA	$(I + \gamma J_o)v$
SLA	$(I + \gamma J)v$
LOLA	$(I + \gamma J_o)v - \gamma \text{diag } J_o^\top \nabla u$
LSS	$(I + J^\top J^{-1})v$
PCGD	$(I - \gamma J_o)^{-1}v$
ED	$-\nabla_x \sup_y \text{NI}(x, y) = -\nabla_x \text{Expl}(x)$
GNI	$-\nabla_x \text{NI}(x, x + \gamma v)$

1105 **Consensus optimization** (CO) (Mescheder, Nowozin, and Geiger, 2017) penalizes the magnitude
 1106 of the simultaneous gradient, encouraging “consensus” between players that attracts them to fixed
 1107 points.

1108 **Symplectic gradient adjustment** (SGA) (Balduzzi et al., 2018), also known as Crossing-the-
 1109 Curl (Gemp and Mahadevan, 2018), reduces the rotational component of game dynamics by using
 1110 the antisymmetric part of the Jacobian.

1111 **Lookahead** (LA) (Zhang and Lesser, 2010) excludes the diagonal components of the Jacobian.
 1112 Each player predicts the behavior of other players after a step of naive learning, but assumes this
 1113 step will occur independently of the current optimization.

1114 In **symmetric lookahead** (SLA) (Letcher, 2018), instead of best-responding to opponents’
 1115 learning, each player responds to *all* players learning, including themselves. It is a linearized version
 1116 of EG (Domingo-Enrich, 2019, Lemma 1.35).

1117 In **learning with opponent-learning awareness** (LOLA) (Foerster et al., 2018), a learner
 1118 optimizes its utility assuming the opponent will take one naive learning step, rather than optimizing
 1119 under the current parameters.

1120 Mazumdar, Sastry, and Jordan (2025) proposed **local symplectic surgery** (LSS) to find local
 1121 NE in two-player zero-sum games. It solves a linear system on each timestep, which is prohibitive
 1122 for high-dimensional parameter spaces. Hence, the authors proposed a two-timescale approximation
 1123 that updates the strategy profile while simultaneously improving an approximate solution to the
 1124 linear system.

1125 **Competitive gradient descent** (CGD) (Schäfer and Anandkumar, 2019) naturally generalizes
 1126 gradient descent to the two-player setting. On each iteration, it jumps to the NE of a quadratically-
 1127 regularized bilinear local approximation of the game. Its convergence and stability properties are
 1128 robust to strong interactions between the players without adapting the stepsize.

1129 **Polymatrix competitive gradient descent** (PCGD) (Ma et al., 2021) generalizes CGD
 1130 to more than two players by jumping to the NE of a quadratically-regularized local polymatrix
 1131 approximation of the game. The series expansion of PCGD to zeroth and first order in γ yields SG
 1132 and LA (Willi et al., 2022, Proposition 4.4), respectively, since $(I - \gamma M)^{-1} = I + \gamma M + \gamma^2 M^2 + \dots$
 1133 for sufficiently small γ . Both CGD and PCGD require solving a linear system of equations on each

1134 iteration, which is prohibitive for high-dimensional parameter spaces (Ma et al., 2021, p. 10).

1135 Lockhart et al. (2019) introduced **exploitability descent** (ED), which directly minimizes
1136 exploitability against worst-case opponents to compute approximate equilibria in two-player zero-
1137 sum extensive-form games. They proved that when both players employ this optimization, the
1138 strategy profile asymptotically converges to an equilibrium, unlike extensive-form fictitious play
1139 (Heinrich, Lanctot, and Silver, 2015) and counterfactual regret minimization (Zinkevich et al., 2007),
1140 where convergence only pertains to time-average strategies. However, ED requires computing exact
1141 best responses y on each iteration, which is generally inefficient or intractable.

1142 To address the intractability of exact best responses, Timbers et al. (2022) introduced approx-
1143 imate exploitability, which utilizes an approximate best response computed through search and
1144 reinforcement learning. This effectively generalizes the local best response metric used in poker (Lisý
1145 and Bowling, 2017).

1146 Similarly, **gradient-based Nikaido–Isoda** (GNI) (Raghunathan, Cherian, and Jha, 2019)
1147 minimizes a local approximation of exploitability using local best responses $y = x + \gamma v$.

1148 Goktas and Greenwald (2022a) recast the exploitability-minimization problem as a min-max
1149 optimization problem and obtained polynomial-time first-order methods for computing variational
1150 equilibria in convex-concave cumulative regret pseudo-games with jointly convex constraints. They
1151 presented two algorithms called **extragradient descent ascent** (EDA) and **augmented descent**
1152 **ascent** (ADA). We benchmark against EDA but not ADA because, unlike the other baselines, ADA
1153 requires *multiple* substeps of gradient ascent *per timestep* to approximate a best response, making
1154 its accuracy highly dependent on the quality of that inner search.

1155 Fiez, Chasnov, and Ratliff (2019) investigated the convergence of learning dynamics in Stackelberg
1156 games with continuous action spaces, characterizing conditions under which attracting critical points
1157 of simultaneous gradient ascent are Stackelberg equilibria in zero-sum games. They developed
1158 a gradient-based update for the leader, paired with a follower employing either a best response
1159 strategy or a gradient-play update rule, yielding algorithms that provably converge to a Stackelberg
1160 equilibrium given appropriate initialization.

1161 Fiez, Chasnov, and Ratliff (2020) studied learning in Stackelberg games, establishing connections
1162 between Nash and Stackelberg equilibria alongside the limit points of simultaneous gradient ascent.
1163 They designed gradient-based learning dynamics that emulate the natural structure of a Stackelberg
1164 game using the implicit function theorem, providing convergence analysis for deterministic and
1165 stochastic updates. However, our focus remains on finding *Nash* equilibria rather than Stackelberg
1166 equilibria.

1167 Fiez, Jin, et al. (2022) considered minimax optimization $\min_x \max_y f(x, y)$ in the context of two-
1168 player zero-sum games, where the min-player minimizes f under the assumption that the max-player
1169 will subsequently maximize it. In their framework, the min-player plays against **smooth algorithms**
1170 deployed by the max-player rather than exact full maximization, which is generally NP-hard. Their
1171 algorithm guarantees monotonic progress, avoids limit cycles, and finds an appropriate stationary
1172 point in a polynomial number of iterations.

1173 Wellman, Tuyls, and Greenwald (2025) surveyed **empirical game-theoretic analysis** (EGTA),
1174 which uses simulation to generate data from which one can induce an empirical game model. The
1175 machine learning approach to empirical game modeling is a form of regression in which the input
1176 is a set of (profile, payoff-vector) pairs and the output is the vector of empirical utility functions.
1177 These techniques can be used to infer a complete empirical game model from an incomplete one.

1178 In a **generative adversarial network** (GAN) (Goodfellow et al., 2014), a generator learns to
1179 produce fake data while a discriminator learns to distinguish it from real data. To stabilize GANs

1180 and address mode collapse, Metz et al. (2017) introduced a method defining the generator objective
1181 with respect to an unrolled optimization of the discriminator. Grnarova, Levy, et al. (2019) proposed
1182 using an approximation of the game-theoretic **duality gap** as a performance measure for GANs, and
1183 later proposed using this measure as the training objective itself to provide convergence guarantees
1184 (Grnarova, Kilcher, et al., 2021).

1185 Gemp, Savani, et al. (2022) proposed **average deviation incentive descent with adaptive**
1186 **sampling** (ADIDAS), which iteratively improves an approximation to a NE through joint play
1187 by tracing a homotopy path defining a continuum of regularized equilibria. To encourage iterates
1188 to remain near this path, the algorithm minimizes the average deviation incentive via stochastic
1189 gradient descent.

1190 Mazumdar, Ratliff, and Sastry (2020) analyzed the limiting behavior of competitive gradient-
1191 based learning algorithms using dynamical systems theory. They characterized a non-negligible
1192 subset of local NE that are actively avoided when agents employ gradient-based learning algorithms.

1193 Mertikopoulos and Zhou (2019) examined the convergence of no-regret learning in continuous
1194 games, focusing on “dual averaging,” where players take small steps along their individual utility
1195 gradients and mirror the output back to their action sets. They introduced the notion of variational
1196 stability, showing that stable equilibria are locally attracting with high probability and globally
1197 stable equilibria are globally attracting with probability 1.

1198 Mertikopoulos and Staudigl (2018) extended this analysis to noisy gradient input, establishing
1199 almost-sure convergence of continuous-time stochastic gradient flows to variationally stable equilibria
1200 under diminishing step sizes. This setting has a connection to our use of zeroth-order pseudo-gradient
1201 estimators (Sections 2.8, 5.1, and 5.3.1), where gradient information is replaced by noisy finite-
1202 difference estimates. Bravo, Leslie, and Mertikopoulos (2018) studied N-player concave games under
1203 bandit feedback, in which each player observes only its own scalar utility rather than a gradient.
1204 They proved no-regret bounds and convergence to Nash equilibria under appropriate monotonicity
1205 conditions. This feedback model closely matches the setting of our JPSPG estimator (Section 5.3.1),
1206 which uses only utility evaluations. Conversely, Vlatakis-Gkaragkounis et al. (2020) established a
1207 negative result: in general finite games, no-regret learning dynamics need not converge to mixed
1208 Nash equilibria, even in the time-average sense. This motivates our use of an exploitability-based
1209 objective and learned best responses (Section 5.2) rather than reliance on time-average no-regret
1210 guarantees.

1211 While most previous work on minimax optimization focused on classical notions of equilibria
1212 from simultaneous games, Jin, Netrapalli, and Jordan (2020) proposed a mathematical definition of
1213 local optimality in sequential game settings, which include GANs and adversarial training. Due to
1214 the nonconvex-nonconcave nature of these problems, minimax is generally not equal to maximin,
1215 making the order in which players act crucial.

1216 To address the rotational behavior of ordinary gradient dynamics in these sequential settings,
1217 Wang, Zhang, and Ba (2020) proposed **follow-the-ridge** (FR) for two-player zero-sum sequential
1218 games, an algorithm that provably converges exclusively to local minimax points.

1219 Tsaknakis and Hong (2021) proposed an algorithm for finding the fractional Nash equilibria
1220 (FNEs) of a two-player zero-sum game with non-convex local cost functions and access only to
1221 local stochastic gradients. Their approach reformulates the problem as minimizing the **regularized**
1222 **Nikaido–Isoda** (RNI) function. Unlike our work, their method is restricted to two-player zero-sum
1223 games and relies on nontrivial subroutines, assuming subproblems are solved to a given accuracy
1224 using external methods like projected gradient descent.

1225 Willi et al. (2022) showed that the original formulation of LOLA is inconsistent because it models

1226 other agents as naive learners rather than LOLA agents, a flaw previously suggested as a cause
1227 for LOLA’s failure to preserve stable fixed points. They formalized consistency and demonstrated
1228 that **higher-order LOLA** (HOLA) solves this inconsistency if it converges. They also proposed
1229 **consistent LOLA** (COLA), a method that relies on no more than second-order derivatives to learn
1230 consistent update functions under mutual opponent shaping, functioning successfully even when
1231 HOLA fails to converge.

1232 Perolat et al. (2022) introduced **DeepNash**, an autonomous agent that mastered the imperfect-
1233 information game Stratego from scratch via self-play using a model-free deep reinforcement learning
1234 method. Its key component, the **Regularised Nash dynamics** (R-NaD) algorithm, modifies the
1235 underlying multiagent learning dynamics to converge to an approximate NE instead of cycling
1236 around it.

1237 Similarly, Qin et al. (2022) proposed PORL, a no-regret style reinforcement learning algorithm
1238 for continuous action tasks with a proven last-iterate convergence guarantee.

1239 Bao and Zhang (2022) proposed **double follow-the-ridge** (double-FTR), an algorithm with
1240 local convergence guarantees to differential NE in general-sum two-player differential games.

1241 Chapter 5

1242 Completed work

1243 In this section, we describe the work we have already completed toward our ultimate goal. Before
1244 doing so, we provide a brief summary of each contribution.

- 1245 1. In Martin and Sandholm (2023), we proposed **randomized policy networks**, a method to
1246 compute approximate Nash equilibria in continuous-action games without access to gradient
1247 information, by combining zeroth-order optimization with randomized policy networks that
1248 model observation-dependent mixed strategies, showing it can efficiently find high-quality
1249 approximate equilibria even when pure-strategy equilibria do not exist and gradients are not
1250 available.
- 1251 2. In Martin and Sandholm (2025b), we proposed **ApproxED: Approximate exploitabil-**
1252 **ity descent via learned best responses**. It consists of two methods that minimize an
1253 approximation of exploitability for continuous-action games by training the strategy profile
1254 jointly with either learned best-response functions or ensembles of candidate best responses,
1255 enabling direct descent on an exploitability proxy. Experiments on continuous games (and
1256 GAN training) show these methods find lower-exploitability strategies and outperform prior
1257 baselines.
- 1258 3. In Martin and Sandholm (2025d), we proposed **joint-perturbation simultaneous pseudo-**
1259 **gradient** (JPSPG), a zeroth-order method for finding approximate Nash equilibria in games
1260 with continuous strategy spaces and no gradient access that performs a single joint perturbation
1261 across all players so the number of utility evaluations per iteration is constant (not linear) in the
1262 number of players. The paper shows this yields large wall-time improvements on many-player
1263 problems (including auctions and other settings with expensive or discontinuous payoffs) and
1264 empirically outperforms per-player perturbation baselines.
- 1265 4. In Martin and Sandholm (2025e), we proposed **player-to-strategy networks** (P2SN), a
1266 neural-network representation that maps player identities to strategies so a single model
1267 can represent strategy profiles across countably or uncountably infinite players by leveraging
1268 function approximation. We also proposed **shared-parameter simultaneous gradient**
1269 (SPSG), an algorithm that generalizes simultaneous gradient ascent to train P2SNs toward
1270 approximate Nash equilibria, and show it converges on a variety of infinite-player games
1271 (including games with infinitely many states/actions and discontinuous utilities).

1272 In each of the following subsections, we describe each contribution in greater detail.

1273 5.1 Finding mixed-strategy equilibria of continuous-action 1274 games without gradients using randomized policy networks

1275 To our knowledge, Martin and Sandholm (2023) was the first work to solve general continuous-action
1276 games with unrestricted mixed strategies and without any gradient information. In that work, we
1277 studied the problem of computing an approximate Nash equilibrium of a continuous-action game
1278 without access to gradients. Such game access is common in reinforcement learning settings, where
1279 the environment is typically treated as a black box. To tackle this problem, we applied zeroth-
1280 order optimization techniques that combine smoothed gradient estimators with equilibrium-finding
1281 dynamics. We modeled players’ strategies using artificial neural networks. In particular, we used
1282 randomized policy networks to model mixed strategies. These take noise in addition to an observation
1283 as input and can flexibly represent arbitrary observation-dependent, continuous-action distributions.
1284 Being able to model such mixed strategies is crucial for tackling continuous-action games that lack
1285 pure-strategy equilibria. We evaluated the performance of our method using an approximation of
1286 the Nash convergence metric from game theory, which measures how much players can benefit
1287 from unilaterally changing their strategy. We applied our method to continuous Colonel Blotto
1288 games, single-item and multi-item auctions, and a visibility game. The experiments showed that
1289 our method can quickly find a high-quality approximate equilibrium. Furthermore, they showed
1290 that the dimensionality of the input noise is crucial for performance. In particular, noise of too low
1291 dimension (or no noise, which yields a deterministic policy) results in high exploitability.

1292 5.1.1 Method

1293 Bichler, Fichtl, Heidekrüger, et al. (2021) modeled strategies using neural networks, where each
1294 player’s policy network takes as input a player’s observation and outputs an action. As described in
1295 Section 4.1.4, these policy networks were then trained using NPGA, which uses Gaussian smoothing
1296 and applies simultaneous gradient ascent. Their policy networks can only model pure strategies,
1297 since the output action is deterministic with respect to the input observation.

1298 We also model strategies using neural networks, with one crucial difference: our policy network f_θ
1299 takes as input the player’s observation o *together with noise* z from some **fixed latent distribution**,
1300 such as the standard multivariate Gaussian distribution. Thus the output $a = f_\theta(o, z)$ of the network
1301 is *random* given o . The network can then learn to transform this randomness into a desired action
1302 distribution. This lets us model mixed strategies, which is especially desirable in games that lack
1303 pure-strategy equilibria. Some approaches in the literature use the output of a policy network to
1304 parameterize some parametric distribution on the action space, such as a Gaussian mixture model.
1305 However, taking the randomness *as input* and letting the neural network *reshape* it as desired allows
1306 us to model arbitrary distributions more flexibly.

1307 Figure 5.1 illustrates the high-level structure of a randomized policy network. It takes as input
1308 an observation and random noise, concatenates them, passes the result through a feedforward
1309 neural network, and outputs an action. The **dimensionality of noise** fed into a randomized policy
1310 network is an important hyperparameter. In Section B.2, we review the literature that studies
1311 the relation between input noise dimension and representational power in neural network-based
1312 generative models.

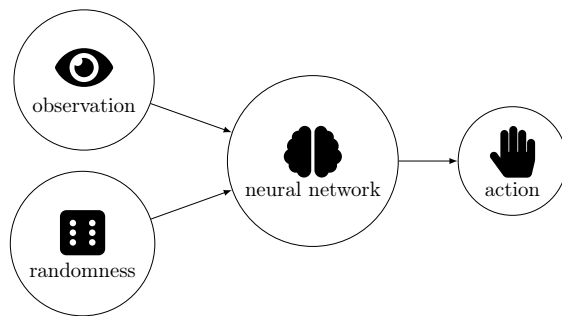


Figure 5.1: High-level structure of a randomized policy network.

1313 Like Bichler, Fichtl, Heidekrüger, et al. (2021), we train these networks through NPGA, which
 1314 can be performed in a massively parallelized or distributed fashion, as described in Algorithm 1.

1315 On any iteration, there is a set of available workers \mathcal{J} . Each worker is assigned the task of
 1316 computing a pseudo-gradient for a particular player. The vector $\{a_j\}_{j \in \mathcal{J}}$ contains the assignment of
 1317 a player for each worker. Each worker’s **pseudorandom number generator (PRNG)** is initialized
 1318 with the same fixed seed. On each iteration, each worker evaluates the utility function (generally the
 1319 most expensive operation and bottleneck for training) twice to compute the finite difference required
 1320 for the pseudo-gradient. It then sends this computed finite difference (a single scalar) to the other
 1321 workers. In turn, it receives the other workers’ scalars. (This is called an “allgather” operation in
 1322 parallel computing.) Therefore, the information that needs to be passed between workers is minimal.
 1323 This greatly reduces the required cross-worker bandwidth compared to schemes that pass parameters
 1324 or gradients between workers, which can be prohibitively expensive for large models.

1325 This massively parallelizes Algorithm 1 of Bichler, Fichtl, Heidekrüger, et al. (2021) (“NPGA using
 1326 ES gradients”). Simultaneously, it generalizes Algorithm 2 of Salimans et al. (2017) (“Parallelized
 1327 Evolution Strategies”), which also uses shared seeds, to the multiplayer setting, with separate gradient
 1328 evaluations and optimizers for each player. Furthermore, it allows for the possibility of setting the
 1329 worker-player assignments a_j and perturbation noise scales ε_j dynamically over time, provided that
 1330 this is done consistently across workers (for example, based on their common state variables). Vanilla
 1331 gradient descent, momentum gradient descent, optimistic gradient descent, or other optimization
 1332 algorithms can be used.

1333 The set of available workers can also change dynamically over time. If a worker leaves or joins
 1334 the pool, the coordinator notifies all workers of its ID so they can remove it from, or add it to, their
 1335 \mathcal{J} sets. The new worker is brought up to speed by passing it the current PRNG state, strategy
 1336 profile parameters, and optimizer states (what state information is needed depends on the algorithm
 1337 used, for example, whether momentum is used).

1338 5.1.2 Experiments

1339 We tested our approach on various benchmark games from the literature, evaluating the exploitability
 1340 of the resulting learned strategy profile. To initialize our networks, we use He initialization (He
 1341 et al., 2015), which is widely used for feedforward networks with ReLU-like activation functions. It
 1342 initializes bias vectors to zero and weight matrices with normally-distributed entries scaled by $\sqrt{2/n}$,
 1343 where n is the layer’s input dimension. We use the ELU activation function (Clevert, Unterthiner,

Algorithm 1 Distributed multiagent pseudo-gradient ascent

Input: \mathcal{I} is the set of players, u is the utility function

initialize PRNG state with fixed seed

$\mathbf{x} \leftarrow$ initial strategy profile

for $i \in \mathcal{I}$ **do**

$S_i \leftarrow \text{init}(\mathbf{x}_i)$

▷ initial state of optimizer i

loop

$\mathcal{J} \leftarrow$ set of available workers

for $j \in \mathcal{J}$ **do**

$a_j \in \mathcal{I}$

▷ set a player; can be set dynamically

$\varepsilon_j \in \mathbb{R}_{>0}$

▷ set a scale; can be set dynamically

$\mathbf{z}_j \sim N(\mathbf{0}, \mathbf{I}_{\dim \mathbf{x}_i})$ where $i = a_j$

$j \leftarrow$ own worker ID

$\delta_j \leftarrow \frac{u(\mathbf{x}[i \rightarrow \mathbf{x}_i + \varepsilon_j \mathbf{z}_j])_i - u(\mathbf{x}[i \rightarrow \mathbf{x}_i - \varepsilon_j \mathbf{z}_j])_i}{2\varepsilon_j}$ where $i = a_j$

send δ_j to coordinator, receive δ from coordinator

for $i \in \mathcal{I}$ **do**

$\mathcal{K} \leftarrow \{j \in \mathcal{J} \mid a_j = i\}$

▷ workers assigned i

$\mathbf{v}_i \leftarrow \frac{1}{|\mathcal{K}|} \sum_{j \in \mathcal{K}} \delta_j \mathbf{z}_j$

▷ i 's pseudo-gradient

$S_i, \mathbf{x}_i \leftarrow \text{step}(S_i, \mathbf{v}_i)$

▷ step optimizer i

1344 and Hochreiter, 2016) for hidden layers. Like Bichler, Fichtl, Heidekrüger, et al. (2021), we use 2
1345 hidden layers with 10 neurons each.

1346 We illustrate the performance of our method by plotting the exploitability across optimization
1347 steps. Noise dimensions 0–4 are shown in blue, orange, green, red, and purple respectively. Each
1348 of these is run for 20 trials. Solid lines indicate means across trials. Bands indicate a confidence
1349 interval for this mean with a confidence level of 0.95. These confidence intervals are computed using
1350 bootstrapping (Efron, 1979), specifically the bias-corrected and accelerated (BCa) method (Efron,
1351 1987).

1352 For the gradient estimator, we use the Gaussian distribution with scale $\sigma = 10^{-2}$, $N = 1$
1353 perturbation vectors, and the central-difference stencil (2 evaluations per step). For the optimizer, we
1354 use standard simultaneous gradient ascent with a learning rate of 10^{-6} . To estimate exploitability (see
1355 the BR computation described in Section 3.3), we use 100 observation samples and 300 state samples
1356 (given each observation). We use a 100-point discretization of the action space for the auctions
1357 and visibility game. For the continuous Colonel Blotto games, we use a 231-point discretization
1358 of the action space. It is obtained by enumerating all partitions of the integer 20 into 3 parts and
1359 renormalizing them to sum to 1.

1360 5.1.2.1 Continuous Colonel Blotto

1361 The original **Colonel Blotto game** was introduced by Borel (1921). It is a two-player zero-sum game
1362 in which two players distribute resources over several battlefields. A battlefield is won by whoever
1363 devotes the most resources to it. A player's payoff is the number of battlefields they win. This game
1364 models many real-world situations of conflict or competition that involve resource allocation, such
1365 as political campaigns, research and development, national security, and systems defense.

1366 Various variants have been studied in the literature. Gross and Wagner (1950) analyzed a
 1367 **continuous** variant in which both players have continuous, possibly unequal budgets. They obtained
 1368 exact solutions for various special cases, including all 2-battlefield cases and all 3-battlefield cases with
 1369 equal budgets. Washburn (2013) generalized to the case where battlefield values are unequal across
 1370 battlefields. Kovenock and Roberson (2021) generalized to the case where battlefield values are also
 1371 unequal across players. Adamo and Matros (2009) studied a variant in which players have incomplete
 1372 information about the other player’s resource budgets. Kovenock and Roberson (2011) studied a
 1373 model where the players are subject to incomplete information about the battlefield valuations.
 1374 Boix-Adserà, Edelman, and Jayanti (2021) analyzed the natural multiplayer generalization of the
 1375 continuous Colonel Blotto game.

1376 The general case with continuous allocations, heterogeneous budgets, heterogeneous battlefield
 1377 values across both players and battlefields, and several players is as follows. Suppose there are J
 1378 battlefields. Let b_i be the budget of player i . Let v_{ij} be the value to player i of battlefield j . Player i ’s
 1379 action space is the standard J -simplex dilated by their budget: $A_i = \{a_{ij} \in \mathbb{R} \mid a_{ij} \geq 0, \sum_j a_{ij} = b_i\}$.
 1380 Player i ’s reward function is $r_i(a) = \sum_j v_{ij} w_{ij}(a)$ where w_{ij} is the probability that i wins j . Ties
 1381 are broken uniformly at random.

1382 We test on the 2-player, 3-item continuous Colonel Blotto game with fixed homogeneous budgets
 1383 and valuations. This game was studied by Gross and Wagner (1950), who derived an exact solution.
 1384 They give the following geometric description of the equilibrium strategy: “[The player] inscribes a
 1385 circle within [the triangle] and erects a hemisphere upon this circle. He next chooses a point from a
 1386 density uniformly distributed over the surface of the hemisphere and projects this point straight
 1387 down into the plane of the triangle... He then divides his forces in respective proportion to the
 1388 triangular areas subtended by [this point] and the sides.” The exact solution is illustrated in Figure
 1389 5.2.

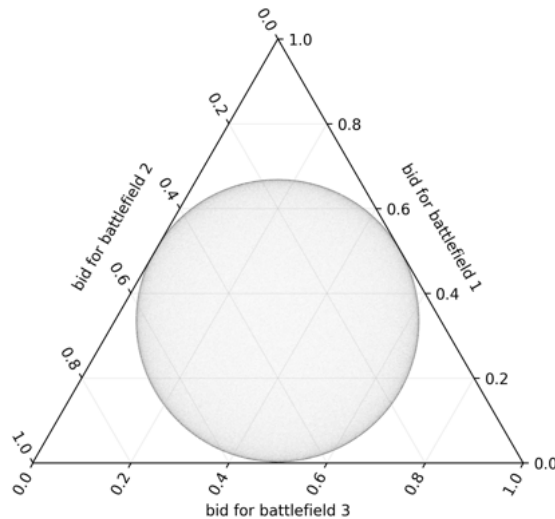


Figure 5.2: Exact solution for 2-player, 3-item continuous Colonel Blotto, shown as a 2D histogram of samples. Darker means higher density.

1390 Actions in the continuous Colonel Blotto game are points on the standard simplex. Thus we use

1391 a softmax activation function for the output layer of the randomized policy network. We describe an
 1392 efficient way to compute best responses in Section B.3.

1393 Figure 5.3 illustrates the performance of our method on the continuous Colonel Blotto game with
 1394 2 players and 3 battlefields. Since the game has no pure-strategy Nash equilibrium, deterministic
 1395 strategies perform badly, as expected. 1-dimensional noise results in slightly better performance, but
 1396 does not let players randomize well enough to approximate the equilibrium. On the other hand, noise
 1397 of dimension 2 and higher is sufficient for good performance. The very slight increase in exploitability
 1398 after 10^8 steps is most likely due to fluctuations introduced by the many sources of stochasticity
 1399 in the training process, including the game and gradient estimates, as well as the fact that we are
 1400 training a multi-layer neural network. Even in the supervised learning setting, loss does not always
 1401 decrease monotonically.

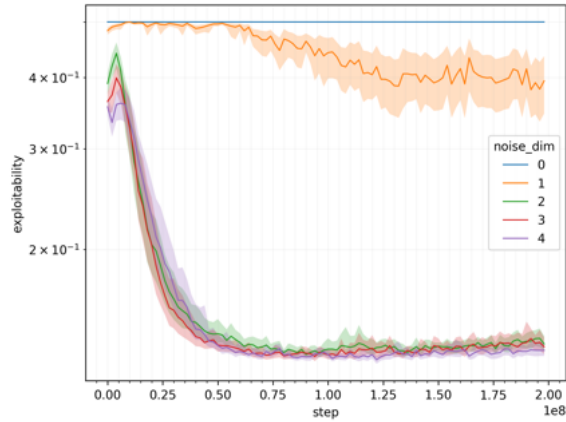


Figure 5.3: Exploitabilities for 2-player, 3-item continuous Colonel Blotto.

1402 Figure 5.4 illustrates the strategies at different stages of training for one trial that uses 2-
 1403 dimensional noise. Each scatter plot is made by sampling 10^5 actions from each player’s strategy.

1404 Figure 5.5 also illustrates performances on the continuous Colonel Blotto game with 2 players and
 1405 3 battlefields. This time, however, the budgets for each player are sampled from the standard uniform
 1406 distribution and revealed to both players. Thus each player must adjust their action distribution
 1407 accordingly. To our knowledge, prior approaches (Adam et al., 2021; Kroupa and Votroubek, 2023;
 1408 Ganzfried, 2021) did not learn strategies that can generalize across different parameters (like budgets
 1409 and valuations), which requires the use of function approximators such as neural networks.

1410 5.1.2.2 Complete-information auction

1411 An auction is a mechanism by which a set of **items** are sold to a set of **bidders**, who have valuations
 1412 for items or sets thereof. Auctions play a central role in the study of markets and are used in a
 1413 wide range of real-world contexts (Krishna, 2002), such as advertising, commodities, radio spectrum
 1414 allocation, real estate, and more. To evaluate their method, Bichler, Fichtl, Heidekrüger, et al. (2021)
 1415 used auctions as a benchmark.

1416 In a **single-item sealed bid auction**, bidders simultaneously submit bids and the highest
 1417 bidder wins the item. Let $w_i(a)$ be the probability i wins given action profile a , where ties are broken

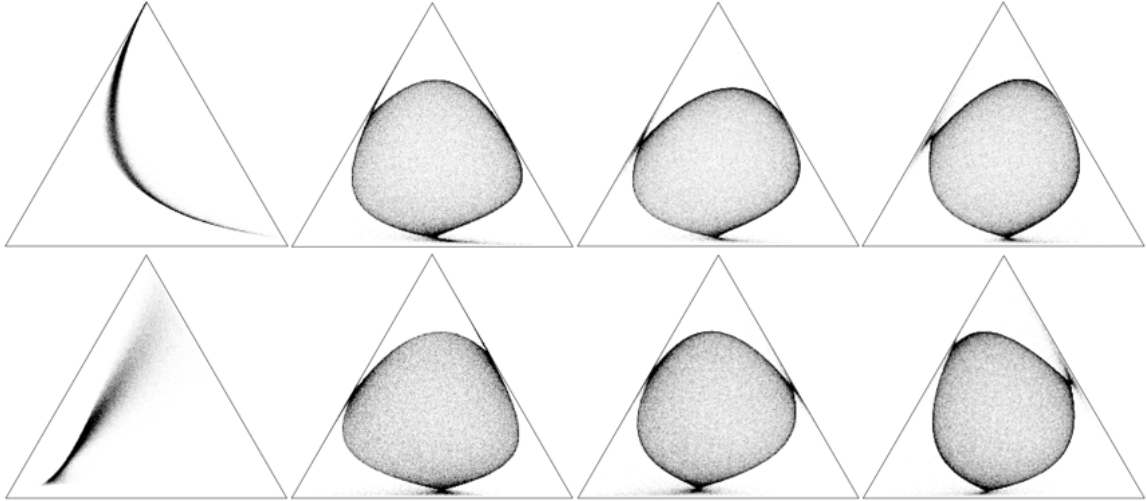


Figure 5.4: Learned strategies for 2-player, 3-item continuous Colonel Blotto, illustrated as 2D histograms of action samples. Top to bottom: Players 1 and 2. Left to right: Across training. Darker means higher density. Each histogram uses 10^4 samples.

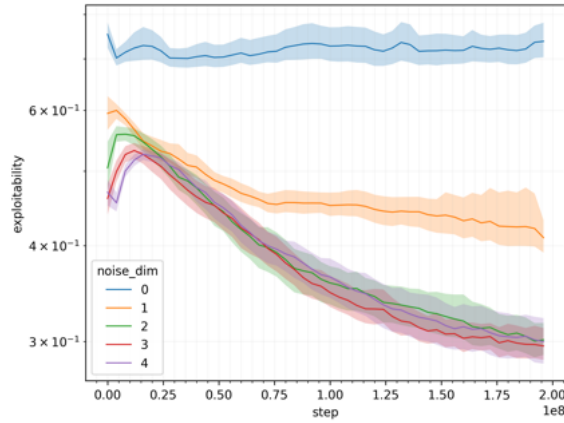


Figure 5.5: Exploitabilities for 2-player, 3-item continuous Colonel Blotto (with random budgets).

1418 uniformly at random. Let $v_i(\omega)$ be the item's value for the i th player given state ω . In a k th-price
 1419 **winner-pay** auction, the winner pays the k th highest bid: $r_i(\omega, a) = w_i(a)(v_i(\omega) - a_{(k)})$, where $a_{(k)}$ is
 1420 the k th highest bid. In an **all-pay** auction, each player always pays their bid: $r_i(\omega, a) = w_i(a)v_i(\omega) - a_i$.
 1421 More details about each type of auction can be found in Section B.1.

1422 The all-pay complete-information auction is widely used to model lobbying for rents in regulated
 1423 and trade protected industries, technological competition and R&D races, political campaigns, job
 1424 promotions, and other contests (Baye, Kovenock, and Vries, 1996). It lacks pure-strategy equilibria

1425 (Baye, Kovenock, and Vries, 1996). We test our method on this auction with 2 players.

1426 This game was studied by Baye, Kovenock, and Vries (1996), who derived an exact solution.
 1427 They show that with homogeneous valuations ($v_1 = v_2 = \dots = v_n$) there exists a unique symmetric
 1428 equilibrium and a continuum of asymmetric equilibria. All of these equilibria are payoff equivalent,
 1429 as is the expected sum of the bids (revenue to the auctioneer). In the symmetric equilibrium, each
 1430 player randomizes uniformly on $[0, v_1]$. The exact solution is illustrated in Figure 5.6.

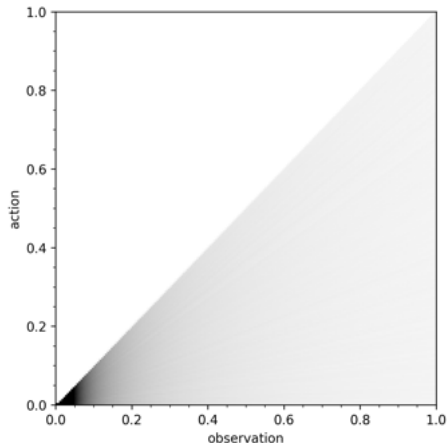


Figure 5.6: Exact solution for the 2-player complete-information auction.

1431 To ensure the output is non-negative, we use a squaring function in the output layer, rather than
 1432 a ReLU function like Bichler, Fichtl, Heidekrüger, et al. (2021). The reason is that, as we found in
 1433 our experiments, ReLU can easily cause degenerate initializations: if the randomly-initialized neural
 1434 network happens to map all of the unit interval (the observation space) to negative bids, no gradient
 1435 signal can be received and the network is stuck.

1436 Figure 5.7 and Figure 5.8 illustrate performances and strategies for the complete-information
 1437 all-pay auction. Recall that these auctions have no pure-strategy equilibria. Thus, as expected,
 1438 deterministic strategies perform poorly. As with Colonel Blotto games, our experiments in these
 1439 auction settings show that the ability to flexibly model mixed strategies is crucial for computing
 1440 approximate Nash equilibria in certain auction settings.

1441 5.1.2.3 Asymmetric-information auction

1442 The 2-player 1st-price winner-pay asymmetric-information auction (in which bidder 1 is perfectly
 1443 informed of the common value of the object, and bidder 2 is completely uninformed) also lacks pure-
 1444 strategy equilibria (Krishna, 2002, Section 8.3). In particular, the second player (who is uninformed)
 1445 must randomize its bid.

1446 This game was studied by Krishna (2002, Section 8.3), who derived an exact solution. Bidder
 1447 1 bids according to the strategy $\beta(v) = E[V \mid V \leq v]$. In our case, $V \sim \mathcal{U}([0, 1])$, so $\beta(v) = \frac{v}{2}$.
 1448 Bidder 2 chooses a bid at random from the interval $[0, E[V]]$ according to the distribution defined
 1449 by $H(b) = P[\beta(V) \leq b]$. In our case, this reduces to the distribution $\mathcal{U}([0, \frac{1}{2}])$. The exact solution is
 1450 illustrated in Figure 5.9.

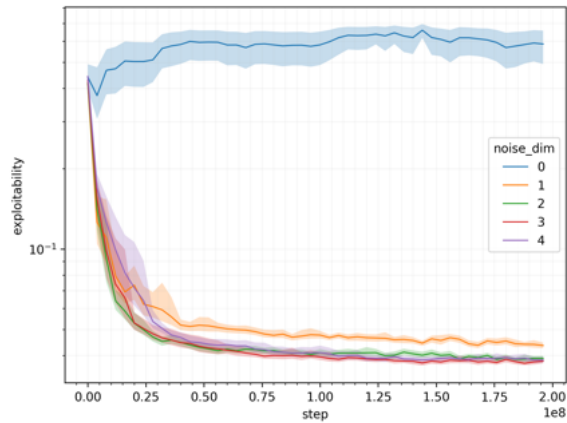


Figure 5.7: Exploitabilities for the 2-player all-pay complete-information auction.

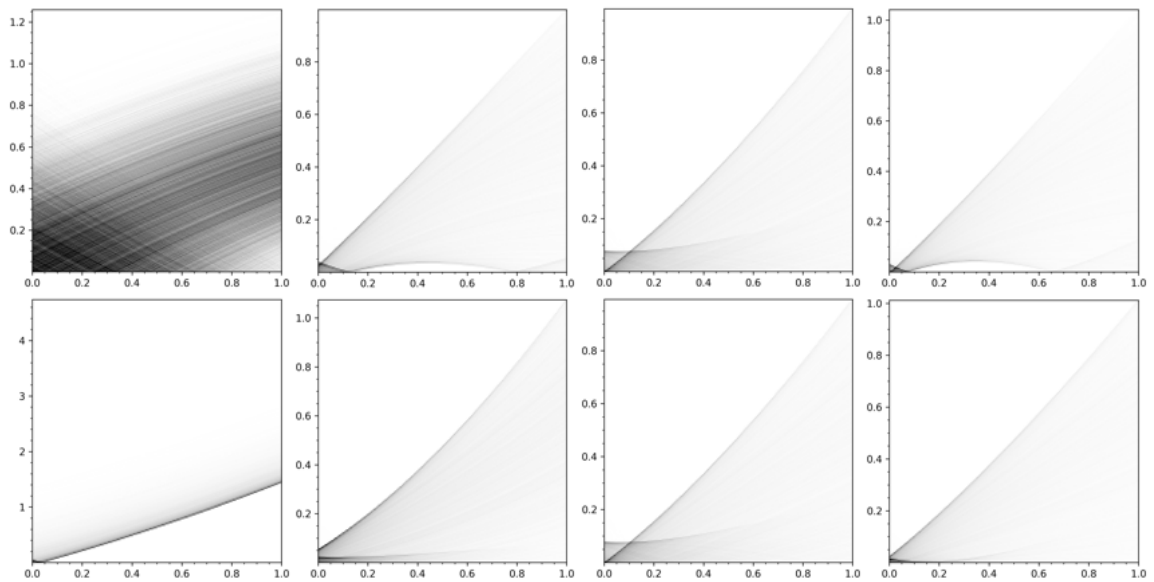


Figure 5.8: Learned strategies for the 2-player all-pay complete-information auction, illustrated as 1D histograms of action samples. Top to bottom: Players 1 and 2. Left to right: Across training. Darker means higher density. X and Y axes denote observation and bid, respectively. Each histogram uses 10^4 samples per observation.

1451

Figure 5.10 illustrates performances for the asymmetric information auction.

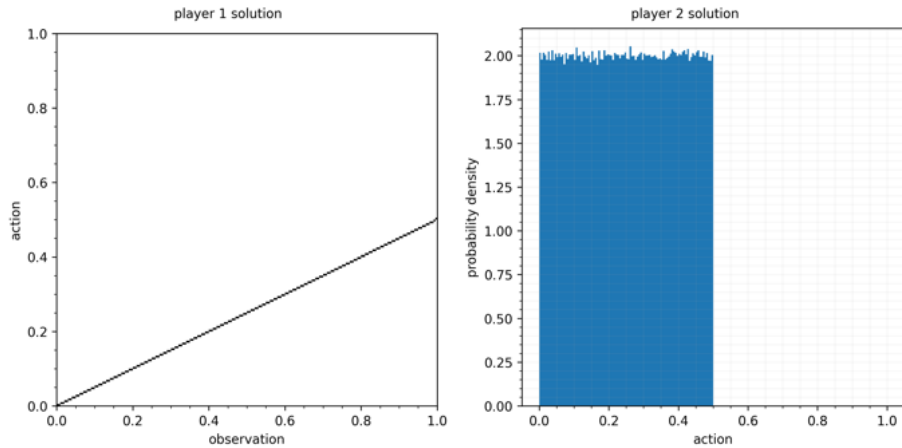


Figure 5.9: Exact solution for the asymmetric-information auction. Left to right: Players 1 and 2.

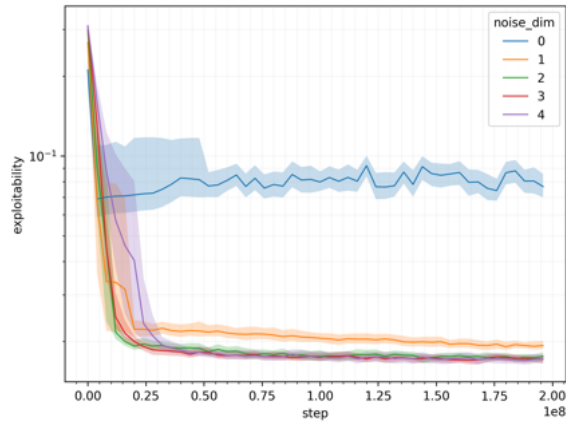


Figure 5.10: Exploitabilities for the asymmetric-information auction.

1452 5.1.2.4 Chopstick auction

1453 Multi-item auctions are of great importance in practice, for example in strategic sourcing (Sandholm,
 1454 2013) and radio spectrum allocation (Milgrom and Segal, 2014; Milgrom and Segal, 2020). However,
 1455 deriving equilibrium bidding strategies for multi-item auctions is notoriously elusive. A rare notable
 1456 instance where equilibrium strategies have been derived is the **chopstick auction** (Szentes and
 1457 Rosenthal, 2003b; Szentes and Rosenthal, 2003a). In this auction, 3 chopsticks are sold simultaneously
 1458 in separate first-price sealed-bid auctions. There are 2 bidders, and it is common knowledge that a
 1459 pair of chopsticks is worth \$1, a single chopstick is worth nothing by itself, and 3 chopsticks are
 1460 worth the same as 2. Here, pure strategies are triples of non-negative real numbers (bids).

1461 This game was studied by Szentes and Rosenthal (2003b) and Szentes and Rosenthal (2003a),
 1462 who derived an exact solution. They describe it as follows: “The supports of the mixtures that

1463 generate the symmetric equilibria in both the first- and second-price cases, turn out to be the
 1464 surfaces of regular tetrahedra, and the distributions themselves turn out to be uniform on these
 1465 surfaces. In addition, in each case all the points inside the tetrahedron are pure best responses to the
 1466 equilibrium mixture.” Mathematically, define the tetrahedron T as the convex hull of the four points
 1467 $(\frac{1}{2}, \frac{1}{2}, 0)$, $(\frac{1}{2}, 0, \frac{1}{2})$, $(0, \frac{1}{2}, \frac{1}{2})$, and $(0, 0, 0)$. Then the uniform probability measure on the 2-dimensional
 1468 surface of T generates a symmetric equilibrium. Furthermore, all points inside the tetrahedron are
 1469 pure best responses to this equilibrium mixture. The exact solution is illustrated in Figure 5.11.

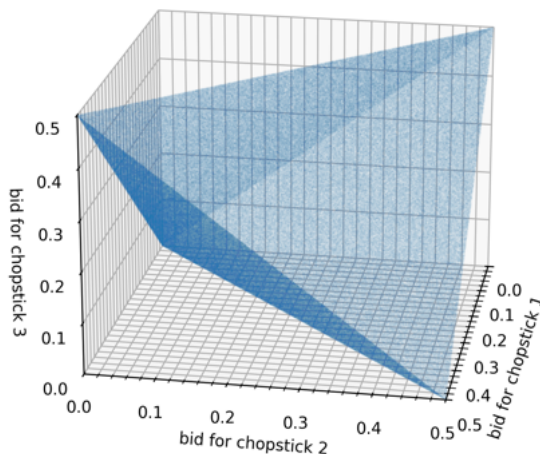


Figure 5.11: Exact solution for the chopstick auction, shown as a 3D scatterplot of samples.

1470 We benchmark on the chopstick auction since it is a rare case of a multi-item auction with a
 1471 known analytic equilibrium, so we can compare our output to an exact equilibrium. It is also a
 1472 canonical case of simultaneous separate auctions under combinatorial preferences.

1473 Figure 5.12 and Figure 5.13 illustrate performances and strategies for the chopstick auction. The
 1474 latter figure shows that, with more epochs, the strategies better approximate a tetrahedron, which
 1475 is the analytic equilibrium (as discussed below).

1476 Here we encounter an interesting phenomenon. Recall that this game has a symmetric equilibrium
 1477 generated by the uniform measure on the surface of a tetrahedron. Although the tetrahedron itself is
 1478 3-dimensional, its surface is only 2-dimensional. Thus one may wonder whether 2-dimensional noise is
 1479 sufficient, that is, whether the network can learn to project this lower-dimensional manifold out into
 1480 the third dimension while “folding” it in the way required to obtain the surface of the tetrahedron.
 1481 Through our experiments, we observe that 2-dimensional noise suffices to (approximately) match the
 1482 performance of higher-dimensional noise. Thus the **intrinsic dimension** of the equilibrium action
 1483 distribution (as opposed to the extrinsic dimension of the ambient space in which it is embedded)
 1484 seems to be the decisive factor.

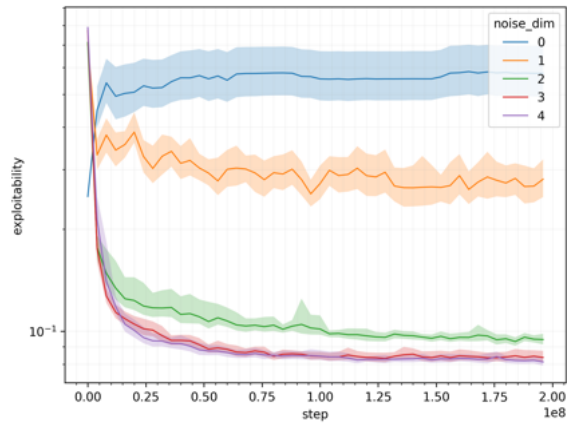


Figure 5.12: Exploitabilities for the chopstick auction.

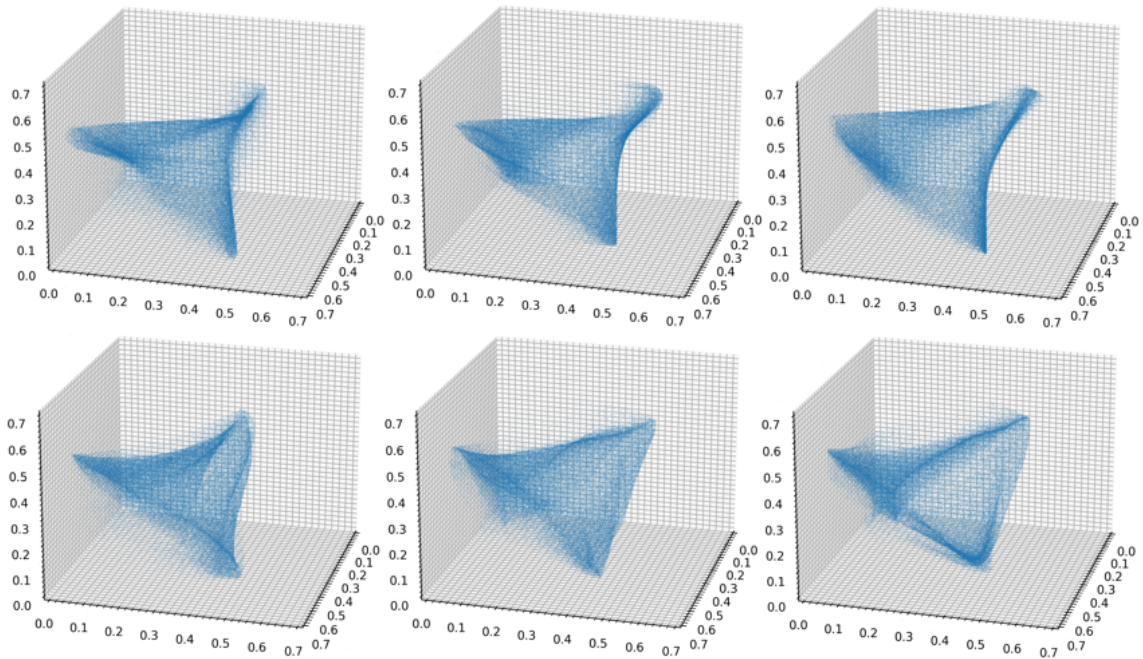


Figure 5.13: Learned strategies for the 2-player, 3-item chopstick auction, shown as 3D scatterplots of action samples. Top to bottom: Players 1 and 2. Left to right: Across training. X, Y, and Z axes denote bid for each item. Each histogram uses 10^5 samples.

1485 **5.1.2.5 Visibility game**

1486 Lotker, Patt-Shamir, and Tuttle (2008) introduced the *visibility game*, a noncooperative, complete-
 1487 information strategic game. In this game, each player i chooses a point $x_i \in [0, 1]$. Their payoff is

1488 the distance to the next higher point, or to 1 if x_i is the highest. This game models a situation
 1489 where players seek to maximize their *visibility time*, and is a variant of the family of “timing games”
 1490 (Fudenberg and Tirole, 1991).

1491 It resembles the “war of attrition” game formalized by Smith (1974). In this game, both players
 1492 are engaged in a costly competition and they need to choose a time to concede. More formally, the
 1493 first player to concede (called “leader”) gets a smaller payoff than the other player (called “follower”).
 1494 Furthermore, the payoff to the leader strictly decreases as time progresses. That is, conceding early
 1495 is better than conceding late.

1496 Lotker, Patt-Shamir, and Tuttle (2008) proved that the n -player visibility game has no pure
 1497 equilibrium, but has a unique mixed equilibrium, which is symmetric. In the 2-player case, up to
 1498 a set of measure zero, there is a unique equilibrium whose strategies have probability densities
 1499 $p(x) = 1/(1 - x)$ when $0 \leq x \leq 1 - 1/e \approx 0.632$ and 0 otherwise. Under this equilibrium, each
 1500 player’s expected payoff is $1/e$. The exact solution is illustrated in Figure 5.14.

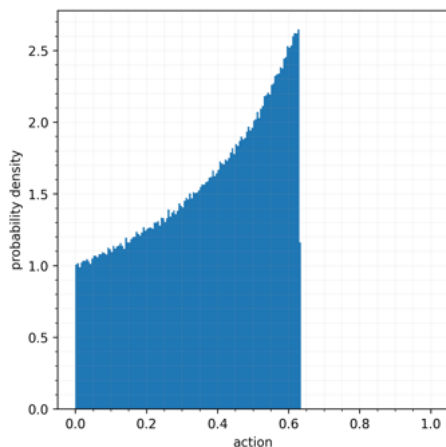


Figure 5.14: Exact solution for the 2-player visibility game.

1501 Figure 5.15 illustrates performances on the 2-player visibility game. Figure 5.16 illustrates
 1502 strategies during training for a trial with 1-dimensional noise. The players’ action distributions
 1503 converge to the expected distribution, including the distinctive cutoff at ≈ 0.632 .

1504 As expected, 0-dimensional noise, which yields deterministic strategies, performs very poorly.
 1505 More interestingly, there is a noticeable gap in performance between 1-dimensional noise, which
 1506 matches the dimensionality of the action space, and higher-dimensional noise. That is, using noise
 1507 of higher dimension than the action space accelerates convergence in this game.

1508 The experiments showed that our method can quickly compute a high-quality approximate
 1509 equilibrium for these games. Furthermore, they showed that the dimensionality of the input noise
 1510 is crucial for representing and converging to equilibrium. In particular, noise of too low dimension
 1511 (or no noise, which yields a deterministic policy) results in failure to converge. Randomized policy
 1512 networks flexibly model observation-dependent action distributions. Thus, in contrast to prior work,
 1513 we can flexibly model mixed strategies and directly optimize them in a “black-box” game with access
 1514 only to payoffs.

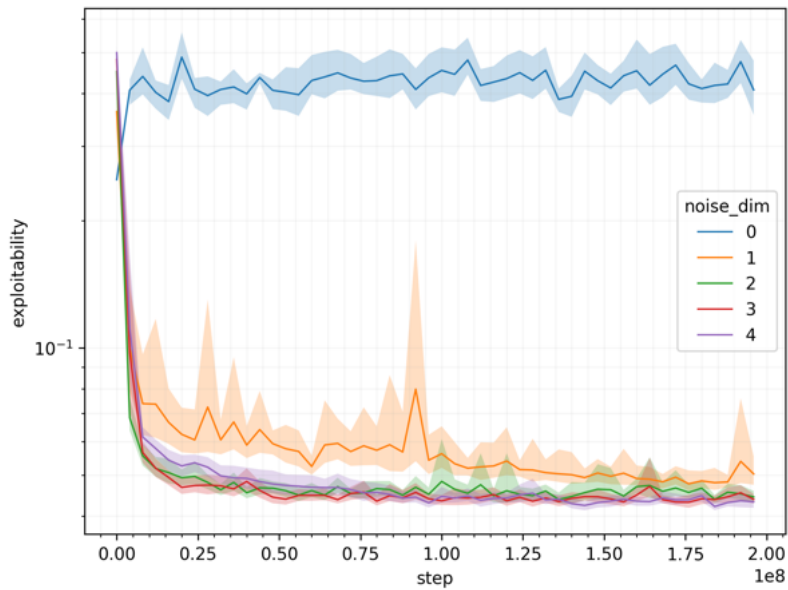


Figure 5.15: Exploitabilities for the 2-player visibility game.

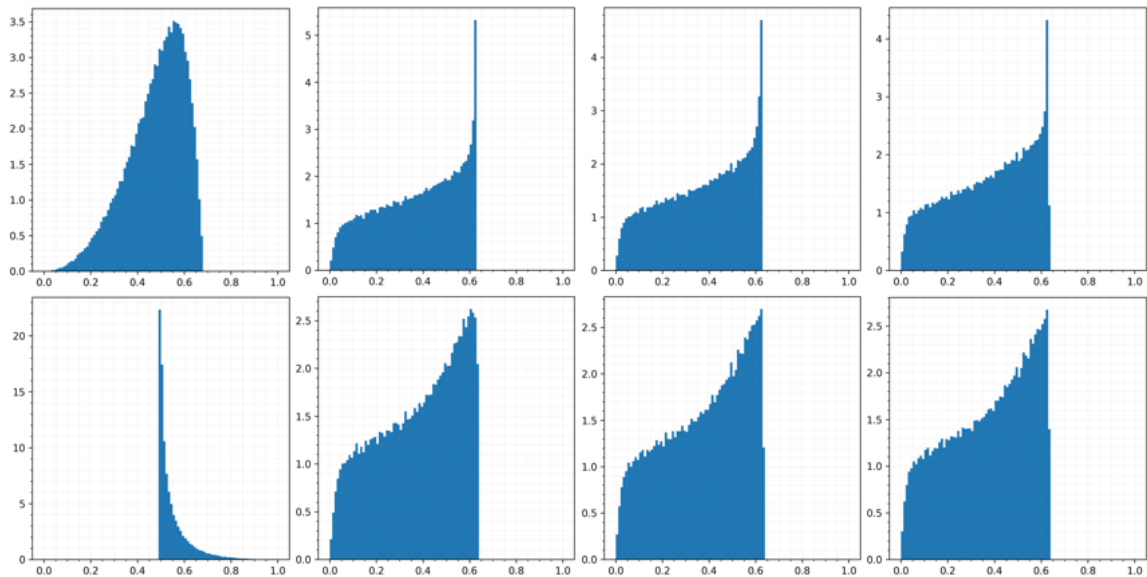


Figure 5.16: Learned strategies for the 2-player visibility game, shown as 1D histograms. Top to bottom: Players 1 and 2. Left to right: Across training. X and Y axes denote the action and probability density, respectively. Each histogram uses 10^5 samples.

1515 5.2 ApproxED: Approximate exploitability descent via learned 1516 best responses

1517 In Martin and Sandholm (2025b), we study the problem of finding an approximate Nash equilibrium
1518 of games with continuous action sets. The standard measure of closeness to Nash equilibrium
1519 is exploitability, which measures how much players can benefit from unilaterally changing their
1520 strategy. We propose two new methods that minimize an approximation of exploitability with respect
1521 to the strategy profile. The first method uses a learned best-response function, which takes the
1522 current strategy profile as input and outputs candidate best responses for each player. The strategy
1523 profile and best-response functions are trained simultaneously, with the former trying to minimize
1524 exploitability while the latter tries to maximize it. The second method maintains an ensemble of
1525 candidate best responses for each player. In each iteration, the best-performing elements of each
1526 ensemble are used to update the current strategy profile. The strategy profile and ensembles are
1527 simultaneously trained to minimize and maximize the approximate exploitability, respectively. We
1528 evaluate our methods on various continuous games and GAN training, showing that they outperform
1529 prior methods.

1530 5.2.1 Method

1531 We are given a utility function u^1 and our goal is to find an NE. Since the exploitability $\text{Expl}(x) =$
1532 $\sup_y \text{NI}(x, y)$ is non-negative everywhere, and zero exactly at NE, we reformulate the problem of
1533 finding an NE (if one exists) as *finding a global minimum of the exploitability function*. That is, we
1534 wish to solve the min-max optimization problem $\inf_x \sup_y \text{NI}(x, y)$. This is equivalent to finding
1535 a minimally-exploitable strategy for a two-player zero-sum meta-game with utility function NI.
1536 To find a minimum, we could try performing gradient descent on NI, like ED does. However, ED
1537 requires best-response oracles. In general games, exact best responses can be difficult or intractable
1538 to compute. Thus we need alternative solutions.

1539 To solve this problem, we can try to perform gradient descent on x and ascent on y simultaneously:
1540 $\dot{x} = -\nabla_x \text{NI}(x, y)$, $\dot{y} = \nabla_y \text{NI}(x, y)$. Unfortunately, this approach can fail even in simple games. For
1541 example, consider the simple bilinear game with $u(x, y) = xy$. The unique Nash equilibrium is at
1542 the origin. However, simultaneous gradient ascent fails to converge to it, and instead cycles around
1543 it indefinitely. The essence of this cycling problem is that Player 2 has to “relearn” a good response
1544 to Player 1 every time Player 1’s strategy switches sign. This is a general problem for games. “Small”
1545 changes in other players’ strategies can cause “large” (discontinuous) changes in a player’s best
1546 response. When such changes occur, players have to “relearn” how to respond to the other players’
1547 strategies. We propose two methods to tackle this problem. These are described in the next two
1548 subsections, respectively.

1549 5.2.1.1 Best-response functions

1550 For this method, we conceptually reformulate the problem as

$$\underset{x \in \mathcal{X}}{\operatorname{argmin}} \sup_{y \in \mathcal{Y}} \text{NI}(x, y) = \underset{x \in \mathcal{X}}{\operatorname{argmin}} \text{NI}(x, b(x)) \quad (5.1)$$

¹For exposition, we assume utility functions are differentiable. If they are not, we can replace gradients with *pseudo-gradients*, as described in Section 2.8.

1551 where $b \in \mathcal{X} \rightarrow \mathcal{Y}$ is a function that satisfies

$$b(x) \in \operatorname{argmax}_{y \in \mathcal{Y}} \text{NI}(x, y) \quad (5.2)$$

1552 Since b is a *function*, it can map different strategies for Player 1 to different strategies for Player 2.
 1553 Thus it can, at least in principle, *immediately* adapt to Player 1’s strategy x , without forgetting
 1554 prior solutions, and thus avoid the cycling problem.

1555 More precisely, suppose \mathcal{Y} is compact and $\text{NI} \in \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is continuous in its second argument.
 1556 Let $x \in \mathcal{X}$. By the extreme value theorem, a continuous real-valued function on a non-empty
 1557 compact set attains its extrema. Therefore, there exists $y \in \mathcal{Y}$ such that $\text{NI}(x, y) = \sup_{y \in \mathcal{Y}} \text{NI}(x, y)$.
 1558 Since this is true for every $x \in \mathcal{X}$, there exists a function $b \in \mathcal{X} \rightarrow \mathcal{Y}$ such that, for every $x \in \mathcal{X}$,
 1559 $\text{NI}(x, b(x)) = \sup_{y \in \mathcal{Y}} \text{NI}(x, y)$. In other words, b is a **best-response function**. Even when \mathcal{Y} is
 1560 not compact and NI does not attain its extrema, one can define a best-response *value* for any
 1561 $x \in \mathcal{X}$ as $\sup_{y \in \mathcal{Y}} \text{NI}(x, y)$, provided the latter exists. In that case, we have the following. Let
 1562 $\varepsilon > 0$ and $x \in \mathcal{X}$. Any function gets arbitrarily close to its supremum (continuity is not required).
 1563 Therefore, there exists a $y \in \mathcal{Y}$ such that $\text{NI}(x, y) + \varepsilon \geq \sup_{y \in \mathcal{Y}} \text{NI}(x, y)$. Therefore, there exists a
 1564 function $b \in \mathcal{X} \rightarrow \mathcal{Y}$ such that, for every $x \in \mathcal{X}$, $\text{NI}(x, b(x)) + \varepsilon \geq \sup_{y \in \mathcal{Y}} \text{NI}(x, y)$. That is b is an
 1565 ε -approximate best-response function.

1566 To find x and b simultaneously, we can perform simultaneous gradient ascent:

$$\dot{x} = -\nabla_x \text{NI}(x, b(x)) \quad (5.3)$$

$$\dot{b} = +\nabla_b \text{NI}(x, b(x)) \quad (5.4)$$

1567 where $\nabla_x \text{NI}(x, b(x))$ is a total (not partial) derivative. That is, the best response function tries to
 1568 *increase* the exploitability while the strategy profile tries to *decrease* it. Since b is a function, Player
 1569 1’s changing behavior poses no fundamental hindrance to it learning good responses and “saving”
 1570 them for later use if Player 1’s behavior changes. It could even learn a good approximation to the
 1571 true best-response function, leaving Player 1 to face a simple standard optimization problem.

1572 If \mathcal{X} is infinite and \mathcal{Y} is nontrivial, $\mathcal{X} \rightarrow \mathcal{Y}$ has infinite dimension. To represent and optimize
 1573 b in practice, we need a finite-dimensional parameterization of (a subset of) this function space.
 1574 In particular, if b is parameterized by θ and is (approximately) surjective onto \mathcal{Y} , then $\text{Expl}(x) =$
 1575 $\sup_{y \in \mathcal{Y}} \text{NI}(x, y) \approx \sup_{\theta \in \Theta} \text{NI}(x, b_\theta(x))$. Inspired by this idea, we propose jointly optimizing x and θ
 1576 according to the following ODE system.

$$\dot{x} = -\nabla_x \text{NI}(x, b_\theta(x)) \quad (5.5)$$

$$\dot{\theta} = +\nabla_\theta \text{NI}(x, b_\theta(x)) \quad (5.6)$$

1577 We call our method **approximate exploitability descent with learned best-response func-**
 1578 **tions** (ApproxED-BRF). Its structure is depicted in Figure 5.17. As the figure illustrates, it maps a
 1579 *strategy profile* to a profile of *best responses* for each player to the other players. It is an all-to-all
 1580 function. We emphasize that this cross-player propagation of information occurs only for the purpose
 1581 of *finding approximate best responses* in order to train *the strategy parameters*, and does not mean,
 1582 for example, that players now get to observe each other’s actions within the real game itself.

1583 The best-response function can take on many possible forms. One possibility is to use a neural
 1584 network. These have the advantages described in Section 1.1. They are also, by far, the most popular
 1585 function approximators used in game solving. Therefore, we use this approach in our experiments.

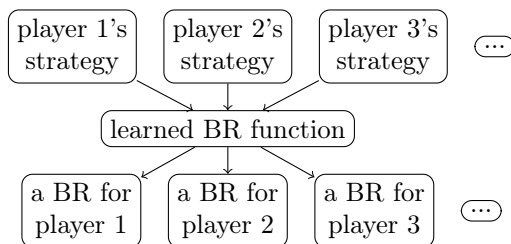


Figure 5.17: Structure of the BRF, in the case of 3 players.

1586 On the other hand, we also experiment with settings where each player’s strategy is *itself*
 1587 represented by a neural network. In this case, the best-response functions take those networks’
 1588 parameters as input. In other words, we adapt the concept of **hypernetworks** (Schmidhuber, 1992;
 1589 Ha, Dai, and Le, 2017; Lorraine and Duvenaud, 2018; MacKay et al., 2019; Bae and Grosse, 2020)
 1590 to the game-theoretic context. We note that, to obtain good performance, the learned best response
 1591 function does not need to represent the true best response function (which may be discontinuous)
 1592 *exactly*, but only *approximately*. Our experimental results indicate that the approximation yielded
 1593 by the neural network performs well across a wide class of games.

1594 5.2.1.2 Best-response ensembles

1595 For this method, we conceptually reformulate the problem as

$$\operatorname{argmin}_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} \text{NI}(x, y) \approx \operatorname{argmin}_{x \in \mathcal{X}} \max_{j \in \mathcal{J}} \text{NI}(x, y_j) \quad (5.7)$$

1596 where \mathcal{J} is a finite set of indices, and $x \in \mathcal{X}$ and $y \in \mathcal{J} \rightarrow \mathcal{Y}$ are trainable parameters. In other
 1597 words, we use an **ensemble** of $|\mathcal{J}|$ responses to x , where the *best* response is selected automatically
 1598 by evaluating x against each y_j and taking the one that attains the best value. Each individual y_j
 1599 is a *strategy* for the original game. Since there are multiple responses in the ensemble, each one
 1600 can “focus on” tackling a particular “type” of behavior from x without having to change drastically
 1601 when the latter changes. We can then train x and y simultaneously: $\dot{x} = -\nabla_x \max_{j \in \mathcal{J}} \text{NI}(x, y_j)$,
 1602 $\dot{y} = \nabla_y \max_{j \in \mathcal{J}} \text{NI}(x, y_j)$. That is, x improves against the best y_j , while the best y_j improves against
 1603 x . Ties are broken in lexicographic order (smaller indices first). This allows for symmetry breaking
 1604 if the ensemble elements are initially equal and the game is deterministic.

1605 There is an issue with the aforementioned scheme, however. If one of the ensemble elements y_j
 1606 strictly dominates the others for all encountered x , then the other elements will never be selected
 1607 under the maximum operator. Thus they will never have a chance to change, improve performance,
 1608 and thus contribute. In that case, the scheme degenerates to ordinary simultaneous gradient ascent.
 1609 We observed this degeneracy in some games.

1610 To solve this issue, we introduced the following approach. To give *all* y_j some chance to
 1611 improve, while incentivizing them to “focus” on particular types of x rather than cover all cases,
 1612 we use a **rank-based weighting** approach. Specifically, we let $\dot{y} = \nabla_y \operatorname{mix}_{j \in \mathcal{J}} \text{NI}(x, y_j)$ where
 1613 $\operatorname{mix}_{j \in \mathcal{J}} a_j = \frac{1}{|\mathcal{J}|} \sum_{j \in \mathcal{J}} r_j a_j$ and $r_j \in \{1, \dots, |\mathcal{J}|\}$ is the ordinal rank of element j . This makes better
 1614 elements receive a higher weight. Thus the best y_j has the most incentive to adapt against the

1615 current x , while others have less incentive, but still some nonetheless.²

1616 Our method is defined by the following ODE system.

$$\dot{x} = -\nabla_x \sum_{i \in \mathcal{I}} \left(\max_{j \in \mathcal{J}} u(x[i \mapsto y_{ij}])_i - u(x)_i \right) \quad (5.8)$$

$$\dot{y} = +\nabla_y \sum_{i \in \mathcal{I}} \left(\text{mix}_{j \in \mathcal{J}} u(x[i \mapsto y_{ij}])_i - u(x)_i \right) \quad (5.9)$$

1617 Here, $y = \{y_{ij}\}_{i \in \mathcal{I}, j \in \mathcal{J}}$ is an ensemble of $|\mathcal{J}|$ responses for each individual player $i \in \mathcal{I}$. We call
 1618 our method **approximate exploitability descent with learned best-response ensembles**
 1619 (ApproxED-BRE). Its structure is depicted in Figure 5.18.

1620 One can compute the $u(x[i \mapsto y_{ij}])_i$ and their gradients in parallel for $i \in \mathcal{I}, j \in \mathcal{J}$. Hence,
 1621 with access to $O(|\mathcal{I}||\mathcal{J}|)$ cores, the method can run approximately as fast as standard simultaneous
 1622 gradient ascent. The ensemble size can be as big as the amount of memory and number of cores or
 1623 workers available allows. If a practitioner has access to parallel computing infrastructure, they can
 1624 scale up the ensemble size as much as possible, until all of the available parallelism is used up, thus
 1625 reaping the benefits of more coverage of the response space while incurring little to no penalty in
 1626 wall time. An interesting direction for future research would be to analyze the effect of the ensemble
 1627 size.³

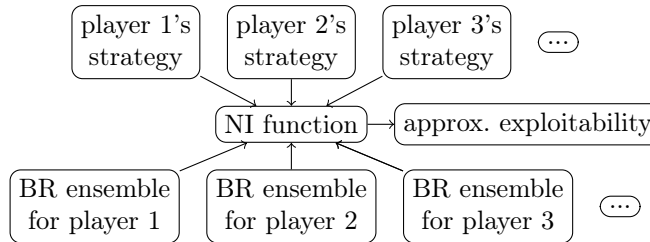


Figure 5.18: Structure of the BRE, in the case of 3 players.

1628 5.2.2 Experiments

1629 In this section, we present some of our experiments. The full list of experiments can be found
 1630 in Martin and Sandholm (2025b). We use a learning rate of $\eta = 10^{-3}$ and $\gamma = 10^{-1}$. For BRF's
 1631 best-response function, we use a fully-connected network with a hidden layer of size 32, the tanh
 1632 activation function, and Glorot weight initialization (Glorot and Bengio, 2010). We do not try to find
 1633 the best neural architecture, because this problem comprises an entire field, may be task-specific, and
 1634 is not the focus of our paper. Thus our experiments are conservative, in the sense that our technique
 1635 could perform even better compared to the baselines if engineering effort were spent tuning the
 1636 neural network. For BRE, we use ensembles of size 10 for each player. For each experiment, we ran

²Furthermore, since the weight of each ensemble element depends only on the rank or order of values, it is invariant under monotone transformations of the utility function.

³For clarity, we emphasize that each individual strategy (including each individual element of an ensemble) could, in general, represent a *mixed* strategy of the original game. Thus the support size of the NE (in terms of *pure strategies* of the original game) is not necessarily directly related to the ensemble size.

1637 64 trials. In our plots, solid lines show the mean across trials, and bands show its standard error.
 1638 For games with stochastic utility functions, we used a batch size of 64.

1639 5.2.2.1 Saddle-point game

1640 For illustration, we start with a simple game. This is a two-player zero-sum game with actions that
 1641 are real numbers and utility function $u_1(x, y) = -u_2(x, y) = xy$. It has a unique NE at the origin.

1642 Our experimental results are shown in Figure 5.19. Our methods converge fastest. Our methods
 also take a more direct path to the equilibrium.

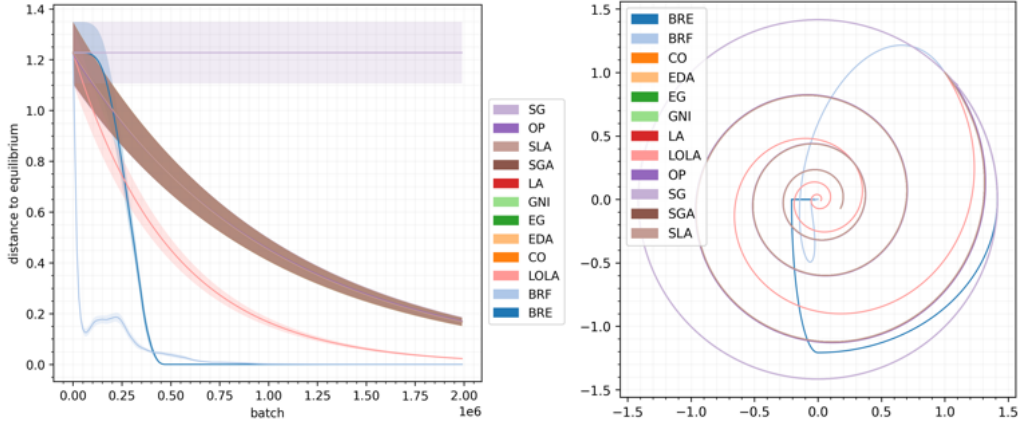


Figure 5.19: Saddle-point game. Left: Distances to equilibrium. Right: Trajectories.

1643

1644 5.2.2.2 Generalized matching pennies

1645 **Matching pennies (MP)** is a classic two-player zero-sum normal-form game with 2 actions per
 1646 player. It can be generalized to n players (Jordan, 1993; Leslie and Collins, 2003) by letting
 1647 $u(a)_i = \llbracket (a_i = a_{i+1 \bmod n}) \oplus (i = n - 1) \rrbracket$, where \oplus denotes exclusive disjunction. That is, each player
 1648 seeks to match the next, but the last player seeks to *unmatch* the first. Utilities are shown in Table
 1649 5.1.

	H	T
H	+1	-1
T	-1	+1

Table 5.1: Utilities for matching pennies (2 players), from the perspective of player 1.

1650 This game has a unique mixed-strategy NE where each player mixes uniformly over its actions. The
 1651 3-player game's NE is locally unstable in a strong sense (Jordan, 1993). More precisely, discrete-time
 1652 fictitious play (Brown, 1951) fails to converge, and instead enters a limit cycle asymptotically.

1653 Our experimental results are shown in Figure 5.20. In the 2-player game, our methods converge
 1654 fastest. In the 3-player game, our methods converge, while the rest diverge.

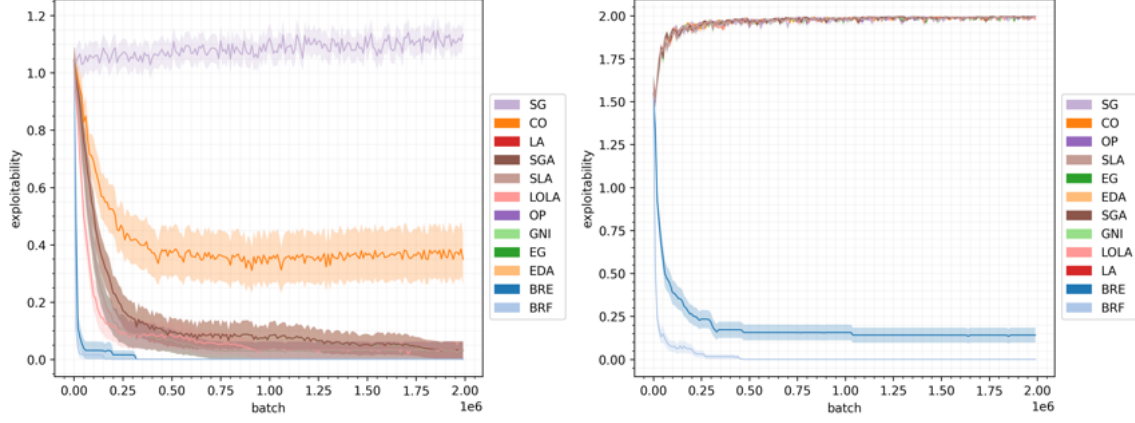


Figure 5.20: Matching pennies with 2 and 3 players.

1655 **5.2.2.3 Generalized rock paper scissors**

1656 **Rock paper scissors** (RPS) is a classic two-player zero-sum normal-form game with 3 actions per
 1657 player. It has a unique mixed-strategy NE where each player mixes uniformly over its actions. Cloud,
 1658 Wang, and Kerr (2023, p. 7) generalize RPS to n actions by letting $u_1(a) = -u_2(a) = \llbracket a_2 - a_1 = 1$
 1659 $\text{mod } n \rrbracket - \llbracket a_1 - a_2 = 1 \text{ mod } n \rrbracket$. Utilities are shown in Table 5.2 and Table 5.3.

	R	P	S
R	0	-1	+1
P	+1	0	-1
S	-1	+1	0

Table 5.2: Utilities for rock paper scissors (3 actions), from the perspective of player 1.

	A	B	C	D
A	0	-1	0	+1
B	+1	0	-1	0
C	0	+1	0	-1
D	-1	0	+1	0

Table 5.3: Utilities for rock paper scissors (4 actions), from the perspective of player 1.

1660 Our experimental results are shown in Figure 5.21. Our methods converge fastest.

1661 **5.2.2.4 Shapley game**

1662 This is a two-player normal-form game with 3 actions per player. Utilities are shown in Table 5.4.

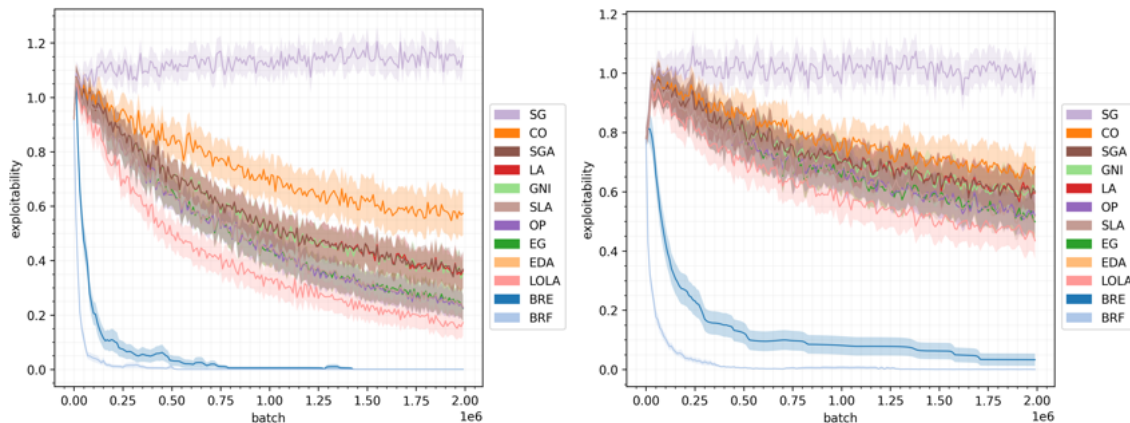


Figure 5.21: Rock paper scissors with 3 and 4 actions.

	A	B	C		A	B	C
A	1	0	0	A	0	1	0
B	0	1	0	B	0	0	1
C	0	0	1	C	1	0	0

Table 5.4: Utilities for the Shapley game (players 1 and 2).

1663 This game was introduced by Shapley (1964, p. 26), and is a classic example of a game for which
 1664 fictitious play (Brown, 1951; Berger, 2007) diverges. (Instead, fictitious play cycles through the cells
 1665 with 1's in them, with ever-increasing lengths of play in each of these cells.)

1666 Our experimental results are shown in Figure 5.22. Our methods converge, while the rest diverge
 1667 and perform almost identically to each other.

1668 5.2.2.5 Glicksberg–Gross game

1669 This is a two-player zero-sum normal-form game with continuous action sets $\mathcal{A}_i = [0, 1]$ and utility
 1670 function $u_1(x, y) = -u_2(x, y) = \frac{(1+x)(1+y)(1-xy)}{(1+xy)^2}$. Glicksberg and Gross (1953) analyzed this game
 1671 and proved that it has a unique mixed-strategy NE where each player's strategy has a cumulative
 1672 distribution function of $F(t) = \frac{4}{\pi} \arctan \sqrt{t}$. To model mixed strategies, we use the following implicit
 1673 density model. We feed a sample from a 1-dimensional standard normal distribution into a fully-
 1674 connected network with one hidden layer of size 32 and output layer of size 1. The output is squeezed
 1675 to the unit interval using the logistic sigmoid function.

1676 Our experimental results are shown in Figure 5.23. Our methods converge fastest.

1677 5.2.2.6 Continuous security game

1678 Security games are used to model defender-adversary interactions in many domains, such as the
 1679 protection of infrastructure like airports, ports, and flights (Kamra, Gupta, Fang, et al., 2018), as

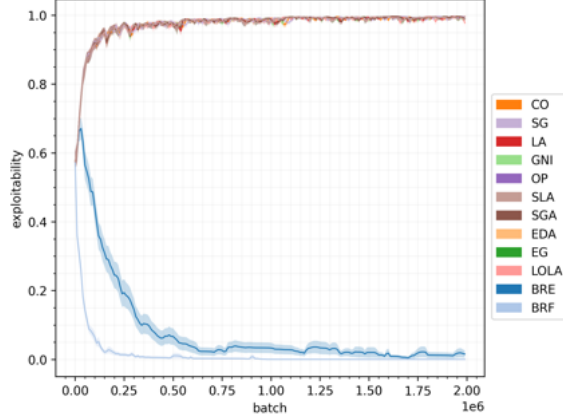


Figure 5.22: Shapley game.

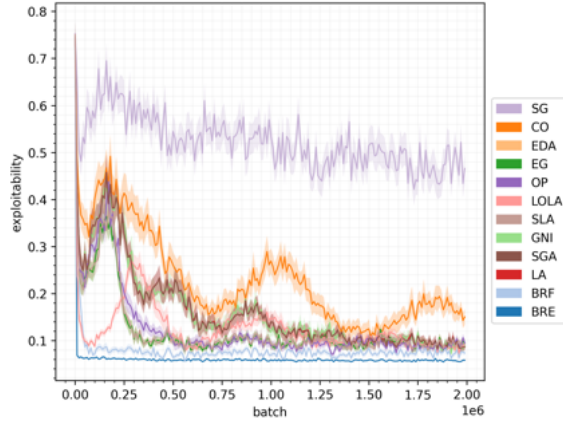


Figure 5.23: Glicksberg-Gross game.

1680 well as wildlife, fisheries, and forests (Kar et al., 2017; Sinha et al., 2018). Security games are often
 1681 modeled with Stackelberg equilibrium as the solution concept, which coincides with NE in zero-sum
 1682 security games and certain structured general-sum games (Korzhyk et al., 2011). Many security
 1683 games have continuous action spaces. These have been studied by Kamra, Fang, et al. (2017), Kamra,
 1684 Gupta, Fang, et al. (2018), and Kamra, Gupta, Wang, et al. (2019). Consider the following game.
 1685 Let $\mathcal{S} = [0, 1]^2$. The attacker chooses a point $x \in \mathcal{S}$. Simultaneously, the defender chooses n points
 1686 $y_i \in \mathcal{S}$. Let $d = \inf_{i \in [n]} \|x - y_i\|$ be the distance between the attacker’s point and the defender’s
 1687 closest point. The defender receives a utility of $\exp(-d^2)$, and the attacker receives $-\exp(-d^2)$.
 1688 Thus the defender seeks to be close to the attacker, while the opposite is true for the attacker.

1689 Our experimental results are shown in Figure 5.24. Our methods perform best.

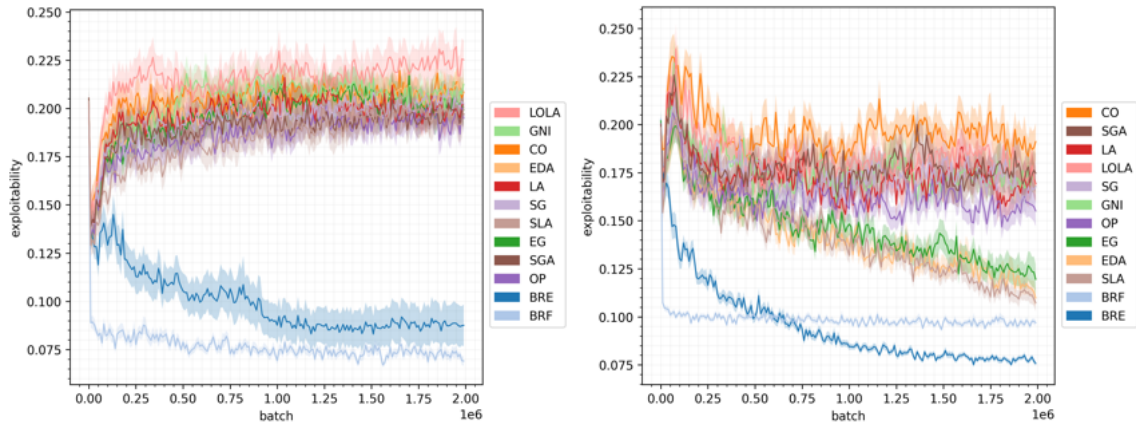


Figure 5.24: Security game with 1 and 2 points.

1690 **5.2.2.7 Poker games**

1691 Kuhn poker is a variant of poker introduced by Kuhn (1950). It is a two-player zero-sum imperfect-
 1692 information game. A 3-player variant was introduced by Szafron, Gibson, and Sturtevant (2013), and
 1693 was one of the largest three-player games to be solved analytically to date. 2-player Kuhn poker has
 1694 a 12-dimensional strategy space per player (24 in total). 3-player Kuhn poker has a 32-dimensional
 1695 strategy space per player (96 in total). Thus the utility function for these games is high-dimensional
 1696 and nonlinear, making them a good benchmark. We use the poker implementation of OpenSpiel
 1697 (Lanctot, Lockhart, et al., 2019).

1698 Our experimental results are shown in Figure 5.25. Our methods converge fastest. The rest
 perform almost identically to each other.

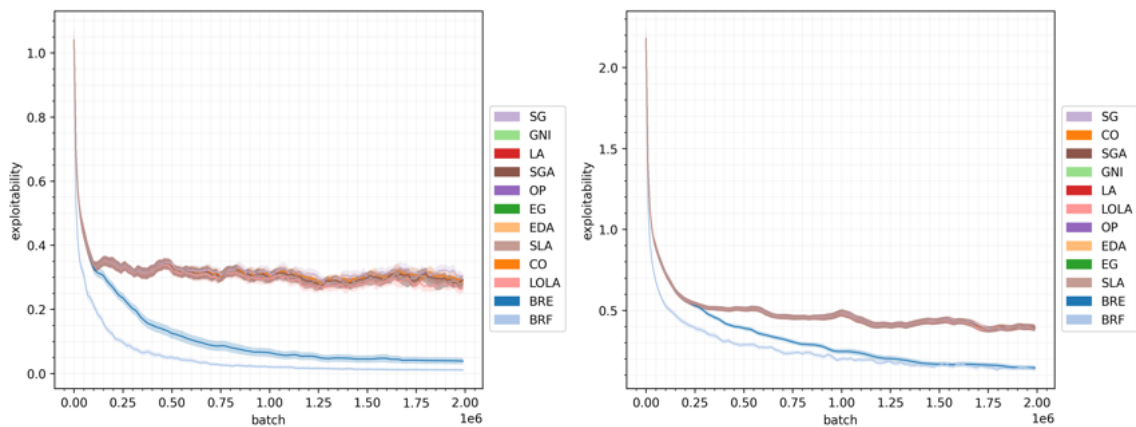


Figure 5.25: Kuhn poker with 2 and 3 players.

1699 **5.2.2.8 GAN training**

1700 A **generative adversarial network** (GAN) (Goodfellow et al., 2014) is a generative model that
1701 consists of two neural networks: a generator and discriminator. The generator maps latent noise
1702 to a data sample. The discriminator maps a data sample to a probability. The generator learns to
1703 generate fake data, while the discriminator learns to distinguish it from real data. GAN training is
1704 a very high-dimensional problem with a highly nonlinear utility function, since the strategies are
1705 parameters for the generator and discriminator.

1706 We test the equilibrium-finding methods on the following datasets. The **ring** dataset consists of
1707 a mixture of 8 Gaussians with a standard deviation of 0.1 whose means are equally spaced around a
1708 circle of radius 1. The **grid** dataset consists of a mixture of 9 Gaussians with a standard deviation
1709 of 0.1 whose means are laid out in a regular square grid spanning from -1 to $+1$ in each coordinate.
1710 The **spiral** dataset consists of a noisy Archimedean spiral, where $t \sim \mathcal{U}(0, 1)$, $r = \sqrt{t}$, $\theta = 2\pi rn$,
1711 $x = \mathcal{N}(r \cos \theta, \sigma)$, and $y = \mathcal{N}(r \sin \theta, \sigma)$. Here, n is the number of turns (we use 2) and σ is the
1712 standard deviation of the noise (we use 0.05). Finally, the **cube** dataset consists of points sampled
1713 uniformly from the edges of a cube and perturbed with Gaussian noise of scale 0.05.

1714 In all cases, the generator’s latent noise distribution is a standard Gaussian matching the
1715 dimension of the dataset. The generator and discriminator have hidden layers of size 32. We
1716 evaluate our method using **Wasserstein distance** (WD) between the real data distribution and
1717 the generator’s fake data distribution. It is the minimum transportation cost needed to turn
1718 one distribution into another. We estimate the WD between the real distribution and generator
1719 distribution by taking 1000 samples of each, computing the Euclidean distance matrix, and solving
1720 the resulting linear assignment problem.⁴

1721 Our experimental results are shown in Figure 5.26. Our methods outperform the rest. In all
1722 cases, SLA diverged to infinity.

1723 We also test on MNIST (Deng, 2012), a dataset of 70,000 28×28 grayscale images of handwritten
1724 digits from 0 to 9. The generator and discriminator networks are the same as before, and fully
1725 connected, but with the hidden layer size increased to 256 and the noise dimension increased to 32.
1726 Due to the larger network size, we use a smaller learning rate of 10^{-4} . Samples are shown in Figure
1727 5.31.

1728 We proposed two new methods that minimize an approximation of the exploitability with respect
1729 to the strategy profile. We evaluated these methods in various continuous games, showing that they
1730 outperform prior methods. In some cases, our methods converge while all prior methods diverge.

⁴For the linear assignment problem, we use the implementation found in the Python scientific computing library SciPy (Virtanen et al., 2020), which uses a modified Jonker-Volgenant algorithm with no initialization (Crouse, 2016).

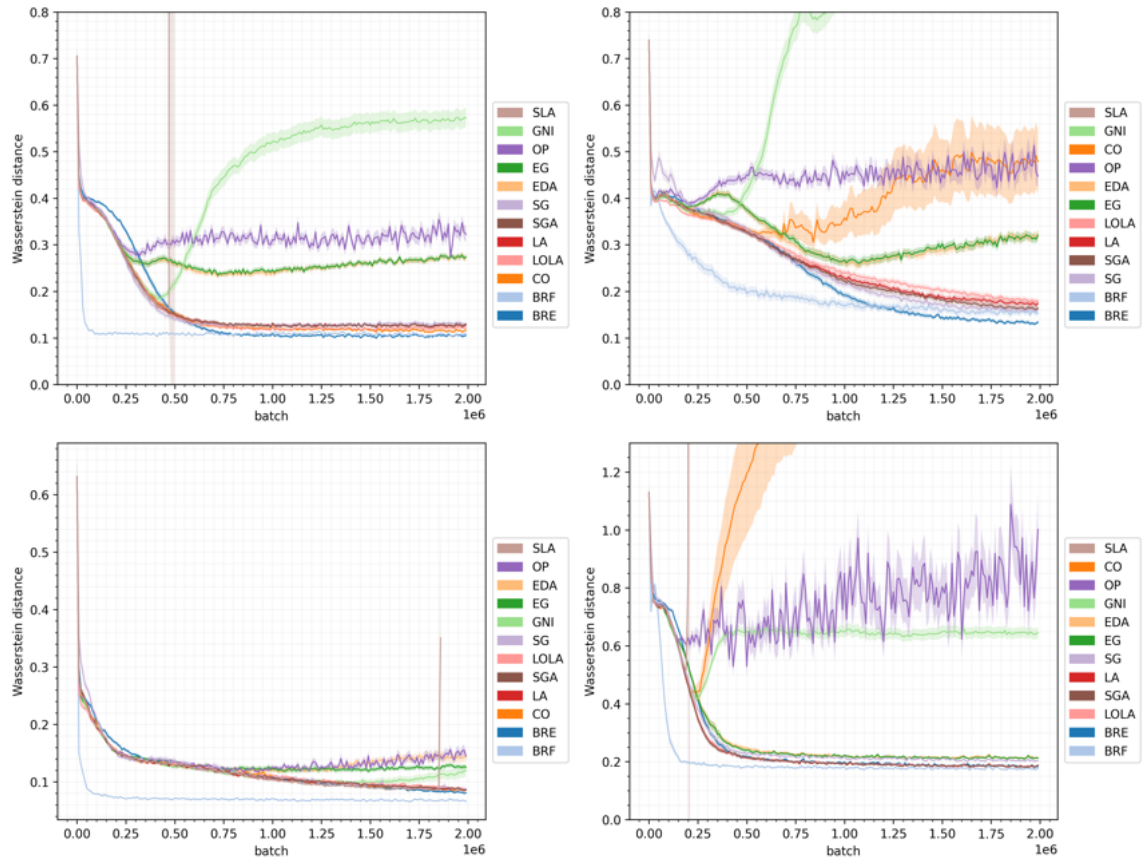


Figure 5.26: Left to right, top to bottom: GAN on ring, grid, spiral, and cube datasets, respectively.

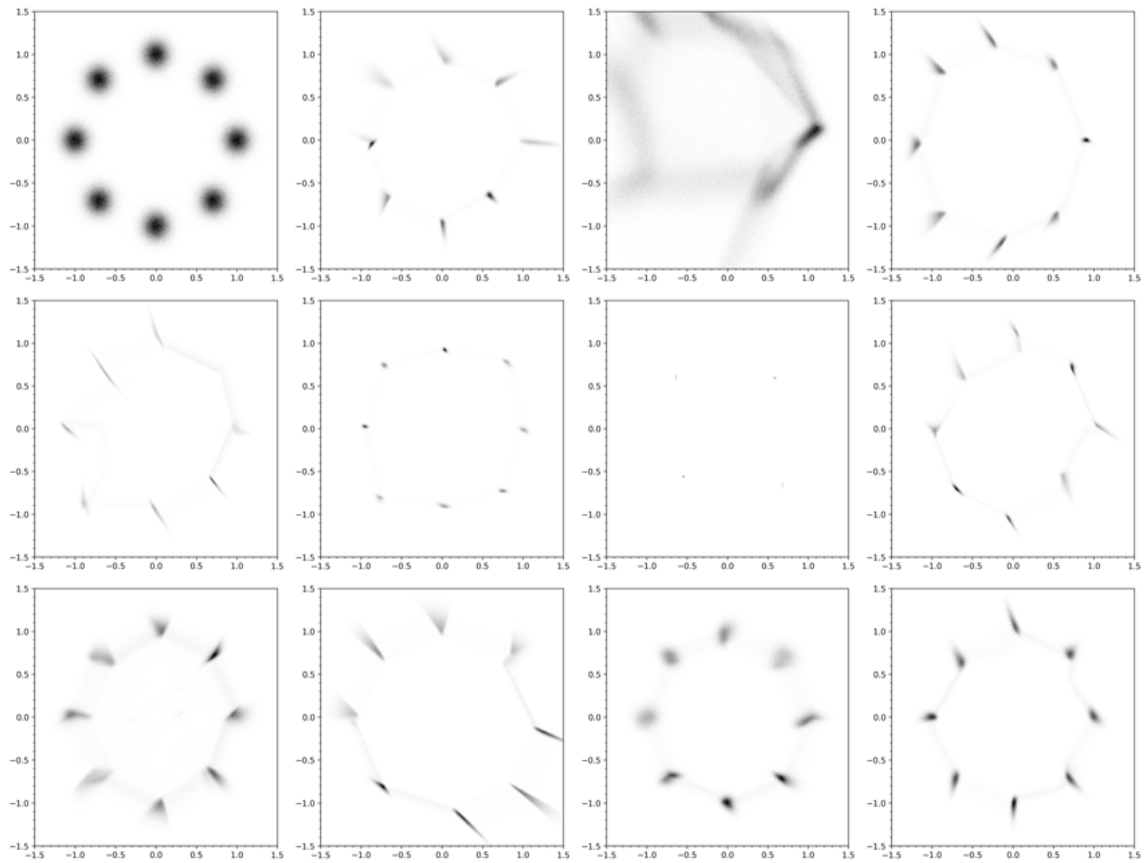


Figure 5.27: GAN (ring dataset). Left to right, top to bottom: Ground truth, SG, OP, EG, CO, SGA, GNI, LA, LOLA, EDA, BRF, BRE. SLA is excluded due to divergence.

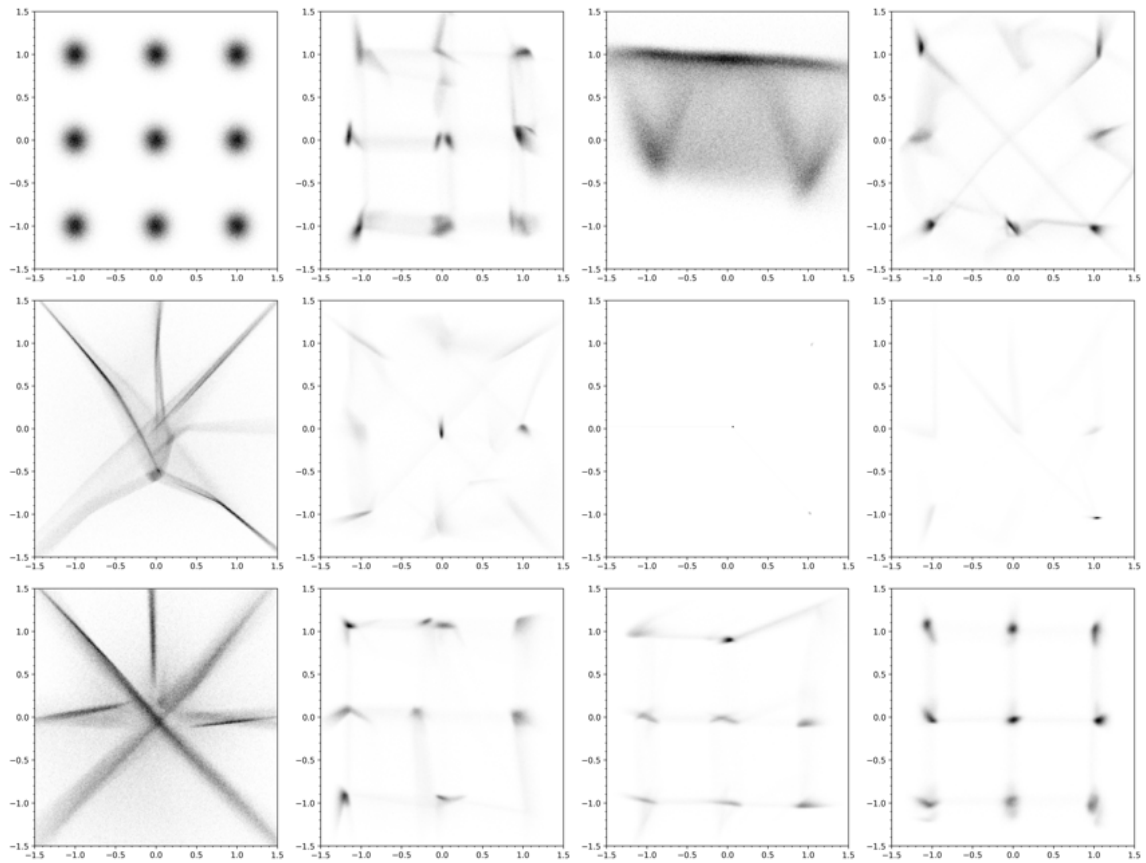


Figure 5.28: GAN (grid dataset). Left to right, top to bottom: Ground truth, SG, OP, EG, CO, SGA, GNI, LA, LOLA, EDA, BRF, BRE. SLA is excluded due to divergence.

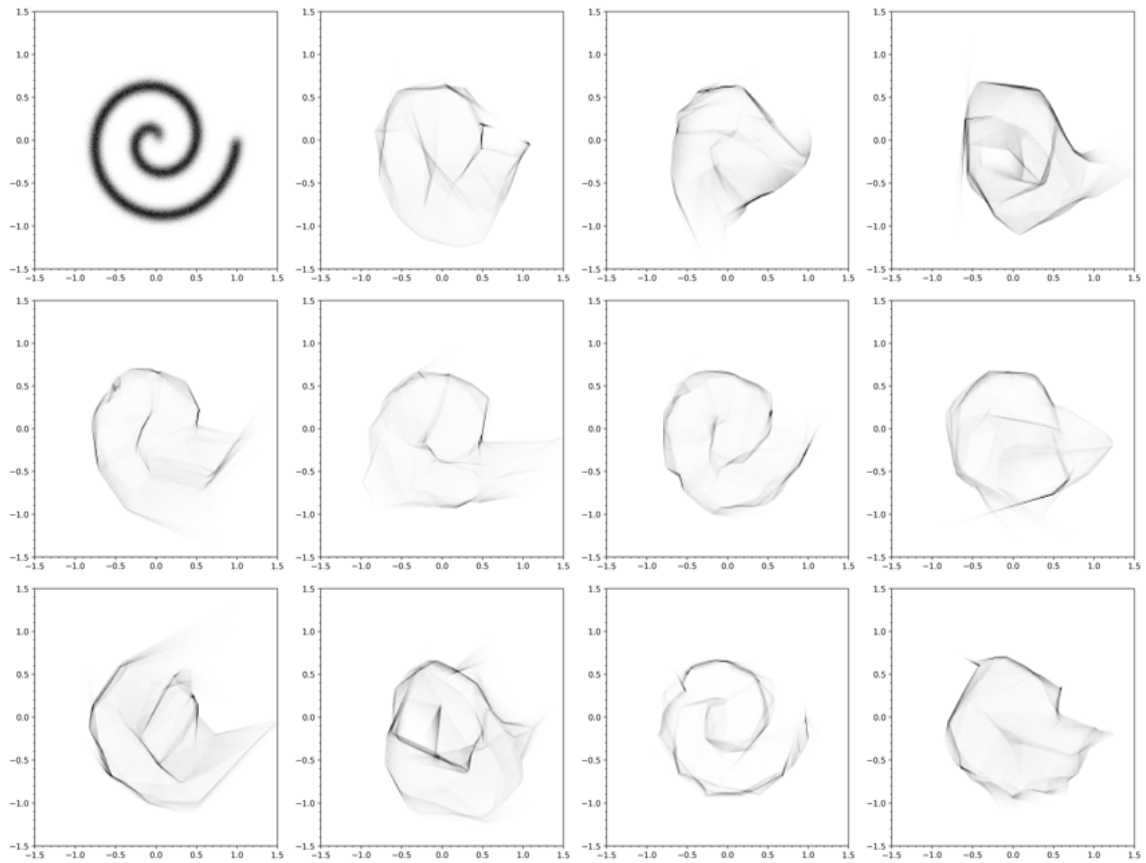


Figure 5.29: GAN (spiral dataset). Left to right, top to bottom: Ground truth, SG, OP, EG, CO, SGA, GNI, LA, LOLA, EDA, BRF, BRE. SLA is excluded due to divergence.

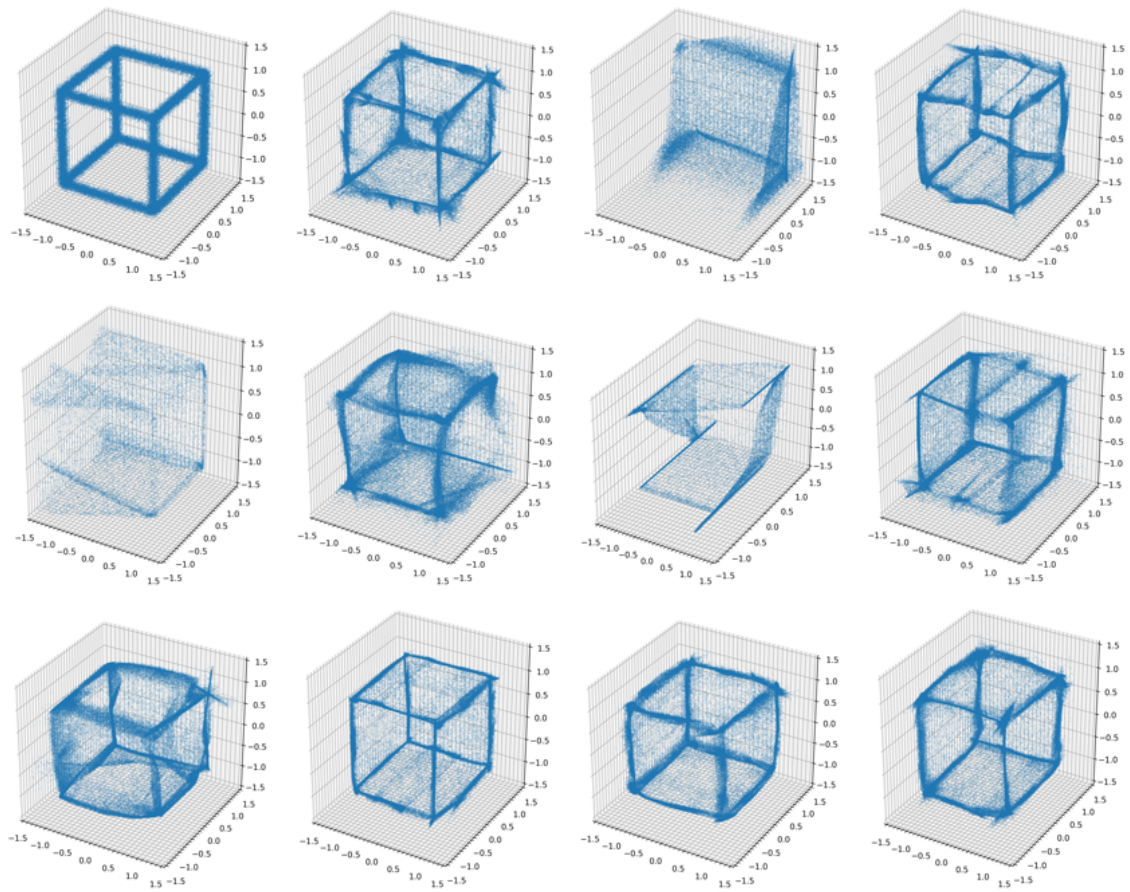


Figure 5.30: GAN (cube dataset). Left to right, top to bottom: Ground truth, SG, OP, EG, CO, SGA, GNI, LA, LOLA, EDA, BRF, BRE. SLA is excluded due to divergence.

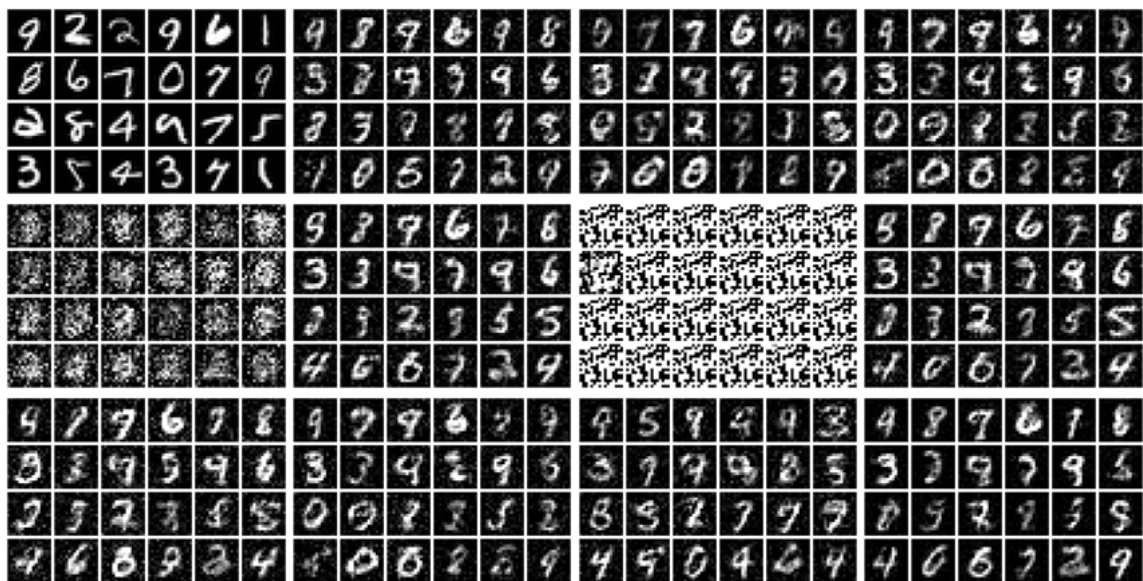


Figure 5.31: GAN (MNIST dataset). Left to right, top to bottom: Ground truth, SG, OP, EG, CO, SGA, GNI, LA, LOLA, EDA, BRF, BRE. SLA is excluded due to divergence.

1731 5.3 Joint-perturbation simultaneous pseudo-gradient

1732 In Martin and Sandholm (2025d), we study the problem of computing an approximate Nash
1733 equilibrium of a game whose strategy space is continuous without access to gradients of the utility
1734 function. Lack of access to gradients is common in reinforcement learning settings, where the
1735 environment is treated as a black box, as well as equilibrium finding in mechanisms such as auctions,
1736 where the mechanism’s payoffs are discontinuous in the players’ actions. To tackle this problem, we
1737 turn to zeroth-order optimization techniques that combine pseudo-gradients with equilibrium-finding
1738 dynamics. Specifically, we introduce a new technique that requires a number of utility function
1739 evaluations per iteration that is constant rather than linear in the number of players. It achieves
1740 this by performing a single joint perturbation on all players’ strategies, rather than perturbing each
1741 one individually. This is very important for many-player games, especially when the utility function
1742 is expensive to compute in terms of wall time, memory, money, or other resources. We evaluate our
1743 approach on various games, including auctions, which have important real-world applications. Our
1744 approach yields a dramatic improvement in performance in terms of the wall time required to reach
1745 an approximate Nash equilibrium.

1746 5.3.1 Method

1747 Recall from Section 4.3 that a common approach to game-solving in the literature is **simultaneous**
1748 **gradient ascent** (SG), which is defined as follows. Let $\mathbf{u} \in \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^n$ be a utility function, where
1749 n is the number of players and d is the dimensionality of each player’s strategy parameters. The
1750 **simultaneous gradient** of \mathbf{u} is the function $\mathbf{v} \in \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ where $\mathbf{v}_i = \nabla_i \mathbf{u}_i$. That is, for each
1751 player, it is the derivative that player’s utility with respect to that player’s parameters. Equivalently,
1752 $\mathbf{v} = \text{diag } \nabla \mathbf{u}$, where $\nabla \mathbf{u}$ is the Jacobian of \mathbf{u} . SG consists of discretizing the ordinary differential
1753 equation $\frac{d}{dt} \mathbf{x} = \mathbf{v}(\mathbf{x})$ in time. That is, each player tries to greedily increase their own utility, acting
1754 as if the other players are fixed. Explicitly, it uses the iteration scheme $\mathbf{x}_{t+1} = \mathbf{x}_t + \alpha_t \mathbf{v}(\mathbf{x}_t)$ for
1755 $t \in \mathbb{N}$, where $\alpha_t > 0$ is some stepsize.

1756 **Simultaneous pseudo-gradient.** SG and other learning dynamics listed in Section 4.3 require
1757 computing the simultaneous gradient \mathbf{v} . However, in some situations, \mathbf{v} does not exist because \mathbf{u} is
1758 not differentiable. In other situations, \mathbf{u} is differentiable, but obtaining an unbiased estimator of
1759 its gradient is difficult or intractable. This can happen if, for example, \mathbf{u} is an expectation (with
1760 respect to a distribution parameterized by \mathbf{x}) of some non-differentiable function. An example of
1761 such a situation is an auction with private values. To resolve this problem, we replace the gradient
1762 of \mathbf{u}_i in the definition of \mathbf{v} with a pseudo-gradient. Explicitly, $\mathbf{g}_i = \frac{1}{\sigma} \mathbf{u}(\mathbf{x}[i \mapsto \mathbf{x}_i + \sigma \mathbf{z}_i])_i \mathbf{z}_i$ for each
1763 Player i , where $\mathbf{z}_i \sim \mu_i$ and μ_i is a multivariate standard normal distribution of the same dimension
1764 as \mathbf{x}_i . That is, we estimate the pseudo-gradient of \mathbf{u}_i (which is a scalar-valued function, since it
1765 outputs only the utility of Player i) with respect to the parameters of Player i . This is the approach
1766 taken by Bichler, Fichtl, Heidekrüger, et al. (2021). It requires one perturbation for each player and
1767 subsequent evaluation of \mathbf{u} . Thus, the number of utility function evaluations per iteration is linear
1768 in the number of players.

1769 **Pseudo-Jacobian.** We extend the preceding concept of the pseudo-gradient from a scalar-valued
1770 function to a vector-valued function. Let $n \in \mathbb{N}$, $\mathbf{f} \in \mathbb{R}^d \rightarrow \mathbb{R}^n$, and $\mathbf{f}_\sigma(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim \mu} \mathbf{f}(\mathbf{x} + \sigma \mathbf{z})$. By

1771 analogy with the pseudo-gradient, we call $\nabla \mathbf{f}_\sigma$ the **pseudo-Jacobian** of \mathbf{f} . Furthermore, it satisfies
 1772 the identity $\nabla \mathbf{f}_\sigma(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim \mu} \frac{1}{\sigma} \mathbf{f}(\mathbf{x} + \sigma \mathbf{z}) \otimes \mathbf{z}$. Therefore, we have an unbiased estimator for it.

1773 In the remainder of this subsection, we show that this estimator is not “too noisy”. (For example,
 1774 it is possible in principle for an estimator to be unbiased, but have very large or even infinite
 1775 variance.) To do this, we give quantitative upper bounds on the moments of the magnitude of this
 1776 estimator. Suppose we use the central-difference stencil described in Section 2.8. This estimator is
 1777 $\mathbf{J} = \frac{\mathbf{f}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{f}(\mathbf{x} - \sigma \mathbf{z})}{2\sigma} \otimes \mathbf{z}$ where \mathbf{z} is a sample from the standard d -dimensional multivariate normal
 1778 distribution. Suppose \mathbf{f} is α -Hölder continuous with constant C . That is, for all $\mathbf{x}, \mathbf{y} \in \text{dom } \mathbf{f}$,
 1779 $\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| \leq C\|\mathbf{x} - \mathbf{y}\|^\alpha$. Then

$$\|\mathbf{J}\|_{\text{F}} = \frac{\|\mathbf{f}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{f}(\mathbf{x} - \sigma \mathbf{z})\|}{2\sigma} \|\mathbf{z}\| \quad (5.10)$$

$$\leq \frac{C\|2\sigma \mathbf{z}\|^\alpha}{2\sigma} \|\mathbf{z}\| \quad (5.11)$$

$$= (2\sigma)^{\alpha-1} C \|\mathbf{z}\|^{\alpha+1} \quad (5.12)$$

1780 where $\|\cdot\|_{\text{F}}$ is the matrix Frobenius norm. Now, the n th raw moment of the chi distribution is
 1781 $\mathbb{E} \|\mathbf{z}\|^n = 2^{n/2} \Gamma(\frac{d+n}{2}) / \Gamma(\frac{d}{2})$, where $d \in \mathbb{N}$ is the dimensionality of \mathbf{z} . Thus

$$\mathbb{E} \|\mathbf{J}\|_{\text{F}}^n \leq (2\sigma)^{n(\alpha-1)} C^n \mathbb{E} \|\mathbf{z}\|^{n(\alpha+1)} \quad (5.13)$$

$$= (2\sigma)^{n(\alpha-1)} C^n 2^{n(\alpha+1)/2} \Gamma(\frac{d+n(\alpha+1)}{2}) / \Gamma(\frac{d}{2}) \quad (5.14)$$

$$= 2^{n(3\alpha-1)/2} \sigma^{n(\alpha-1)} C^n \Gamma(\frac{d+n(\alpha+1)}{2}) / \Gamma(\frac{d}{2}) \quad (5.15)$$

1782 This yields an upper bound on any moment of the norm of the estimator. It can be applied to
 1783 C -Lipschitz functions with $\alpha = 1$, and to C -bounded-range functions with $\alpha = 0$.

1784 **Joint perturbation.** We now combine all of the preceding concepts that have been discussed
 1785 so far. That is, we combine (1) the identity $\mathbf{v} = \text{diag } \nabla \mathbf{u}$, (2) the concept of the pseudo-Jacobian,
 1786 and (3) the identity $\text{diag}(\mathbf{a} \otimes \mathbf{b}) = \mathbf{a} \odot \mathbf{b}$ to obtain an estimator that requires only a *single*, joint
 1787 perturbation across all players. Specifically, let $\mathbf{v}_\sigma = \mathbf{v} * G_\sigma$, that is, the smoothing of \mathbf{v} by a
 1788 Gaussian kernel of width σ . Then⁵

$$\mathbf{v}_\sigma(\mathbf{x}) = (\mathbf{v} * G_\sigma)(\mathbf{x}) \quad (5.16)$$

$$= ((\text{diag } \nabla \mathbf{u}) * G_\sigma)(\mathbf{x}) \quad (5.17)$$

$$= \text{diag}(\nabla \mathbf{u} * G_\sigma)(\mathbf{x}) \quad (5.18)$$

$$= \text{diag } \nabla(\mathbf{u} * G_\sigma)(\mathbf{x}) \quad (5.19)$$

$$= \text{diag } \nabla \mathbf{u}_\sigma(\mathbf{x}) \quad (5.20)$$

$$= \text{diag } \mathbb{E}_{\mathbf{z} \sim \mu} \frac{1}{\sigma} \mathbf{u}(\mathbf{x} + \sigma \mathbf{z}) \otimes \mathbf{z} \quad (5.21)$$

$$= \mathbb{E}_{\mathbf{z} \sim \mu} \frac{1}{\sigma} \text{diag } \mathbf{u}(\mathbf{x} + \sigma \mathbf{z}) \otimes \mathbf{z} \quad (5.22)$$

$$= \mathbb{E}_{\mathbf{z} \sim \mu} \frac{1}{\sigma} \mathbf{u}(\mathbf{x} + \sigma \mathbf{z}) \odot \mathbf{z} \quad (5.23)$$

1789 Therefore, we obtain the following unbiased estimator: $\mathbf{g} = \frac{1}{\sigma} \mathbf{u}(\mathbf{x} + \sigma \mathbf{z}) \odot \mathbf{z}$. In terms of indices, for
 1790 clarity: $\mathbf{g}_i = \frac{1}{\sigma} \mathbf{u}(\mathbf{x} + \sigma \mathbf{z})_i \mathbf{z}_i$. With this new estimator, the number of utility function evaluations per

⁵The interchange of differentiation and integration is justified by the Hölder continuity of \mathbf{u} combined with Gaussian smoothing (Duchi, Bartlett, and Wainwright, 2012).

1791 iteration is now *constant* in the number of players, rather than linear. This dramatically reduces
 1792 the number of utility function evaluations when there are many players. Therefore, it makes game
 1793 solving significantly more efficient in many-player games, especially when the utility function is
 1794 expensive to evaluate in terms of wall time, memory, money, or other resources such as real-world
 1795 experiments. We call our method **joint-perturbation simultaneous pseudo-gradient** (JPSPG),
 1796 in contrast to the classical method, **simultaneous pseudo-gradient** (SPG). To our knowledge,
 1797 this is the first work to define the concept of the pseudo-Jacobian and apply it to the estimation
 1798 of the simultaneous gradient. The intuition behind our method is that we estimate the analogue
 1799 of the gradient (i.e., Jacobian) of the *vector*-valued function that returns utilities for all n players
 1800 *simultaneously* (and then select its diagonal), rather than treating the problem as having n different,
 1801 separate *scalar*-valued functions whose gradients need to be estimated separately. This saves us
 1802 work, because we need just 1 rather than n different function evaluations. For clarity, we emphasize
 1803 that our method does not assume or require symmetry across players. It handles general utility
 1804 functions on general strategy spaces. The utility function does not need to be symmetric across
 1805 players. In fact, the players’ strategy spaces need not be equal. For example, player 1’s parameters
 1806 could be 10-dimensional, player 2’s parameters could be 20-dimensional, and so on. One can use
 1807 our method in combination with any gradient-based learning dynamics from the literature that are
 1808 based on the simultaneous gradient. These include learning dynamics such as standard simultaneous
 1809 gradient ascent, extragradient ascent (Korpelevich, 1976), optimistic gradient ascent (Popov, 1980;
 1810 Daskalakis et al., 2018), and so on. In the games that we tested on, simultaneous gradient ascent
 1811 was sufficient to obtain convergence to an approximate NE. However, we can combine JPSPG with
 1812 any other learning dynamics.

1813 5.3.2 Experiments

1814 We test our approach against the classical approach on several continuous-action games, in particular
 1815 on many kinds of auction and on continuous Goofspiel (which can be thought of as a kind of sequential
 1816 auction with budget constraints). These are many-player games where each utility function evaluation
 1817 is expensive, thus highlighting the problem of interest and showing the benefits of our method.
 1818 Specifically, each utility function evaluation requires solving a linear assignment problem, solving
 1819 an integer linear program, running policies over multiple steps, etc., all of which are expensive
 1820 operations. The games we test on cover various theoretical properties of interest, including various
 1821 dimensionalities for each player’s observation, various dimensionalities for each player’s action,
 1822 single-step vs. multi-step settings, etc. Our experimental hyperparameters are as follows. For each
 1823 experiment, we run 8 trials. In each graph, solid lines show the mean across trials, and bands show
 1824 the standard error of the mean. The classical method is shown in blue, while our method is shown
 1825 in orange. We use a stepsize of 10^{-4} . For the Gaussian smoothing, we use a perturbation scale
 1826 σ of 10^{-1} . We use a batch size per iteration of 256. To update parameters, we use the AdaBelief
 1827 optimizer (Zhuang et al., 2020). For each player’s strategy network, we use a single hidden layer
 1828 of size 64, the ReLU activation function, and He initialization (He et al., 2015) for initializing the
 1829 network’s weights.

1830 5.3.2.1 Multi-item unit-demand auction

1831 Here, we consider a type of multi-item auction called a **unit-demand auction**. In this auction,
 1832 we have n bidders and m non-identical items. Each bidder i has a private valuation $v_{i,j}$ for each
 1833 item j . Furthermore, each bidder has *unit demand*, meaning that its value for a *bundle* of items

1834 is the same as that for the maximum-value item in that bundle: $v_i(S) = \max_{j \in S} v_{ij}$, where S is a
 1835 bundle of items. Housing markets are often given as an example of unit-demand preferences. This
 1836 model was first studied by Shapley and Shubik (1971). This is a special case of a **limited-demand**
 1837 **model** with K units in which each bidder has use for at most $L < K$ units, as described in Krishna
 1838 (2002, §13.4.2 and §13.5.2). The single-unit case corresponds to $L = 1$. For our experiment, we
 1839 use a prior where bidder-item valuations v_{ij} are independently sampled uniformly at random from
 1840 the unit interval. Each player i submits a bid b_{ij} for each item j . To allocate items, our auction
 1841 mechanism assigns items to bidders in a way that maximizes the sum of bids across players. This
 1842 requires solving a **linear assignment problem**, which can be described as follows. Given a bid
 1843 matrix $b \in \mathbb{R}^{n \times m}$, compute a binary assignment $x \in \{0, 1\}^{n \times m}$ that satisfies the following:

$$\text{maximize } \sum_{i \in [n]} \sum_{j \in [m]} b_{ij} x_{ij} \tag{5.24}$$

$$\text{subject to } \sum_{i \in [n]} x_{ij} \leq 1 \quad \forall j \in [m] \tag{5.25}$$

$$\sum_{j \in [m]} x_{ij} \leq 1 \quad \forall i \in [n] \tag{5.26}$$

1844 That is, maximize the sum of values subject to the constraint that each item is assigned to at most
 1845 one bidder and each bidder receives at most one item. An example is illustrated in Figure 5.32.

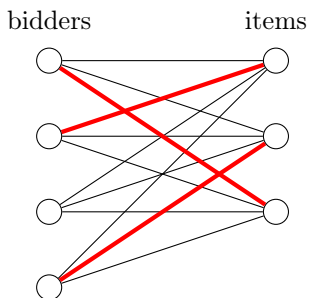


Figure 5.32: Example of a complete bipartite graph between bidders and items, and an assignment.

1846 The linear assignment problem was first described in a seminal paper by Kuhn (1955), who
 1847 introduced a solution approach called the **Hungarian method**. Subsequently, various algorithms
 1848 have been devised in the literature. We use the modified **Jonker-Volgenant algorithm** (Jonker
 1849 and Volgenant, 1988) given in Crouse (2016), as implemented in the Python scientific computing
 1850 library SciPy (Virtanen et al., 2020), and reimplement it in JAX (Bradbury et al., 2018). This
 1851 algorithm has a time complexity of $O(n^3)$.

1852 Figure 5.33 shows the exploitability over the course of training for a unit-demand auction. Both
 1853 our method and the classical method attain a similar exploitability for a given iteration count, but
 1854 our method is dramatically faster in terms of run time (here expressed in seconds). This is explained
 1855 by the fact that the standard method requires many more evaluations of the utility function per
 1856 iteration, each of which requires solving an assignment problem (which is expensive), whereas ours
 1857 only requires a *single* utility function evaluation. The advantage of our method over the baseline
 1858 increases as the number of players increases.

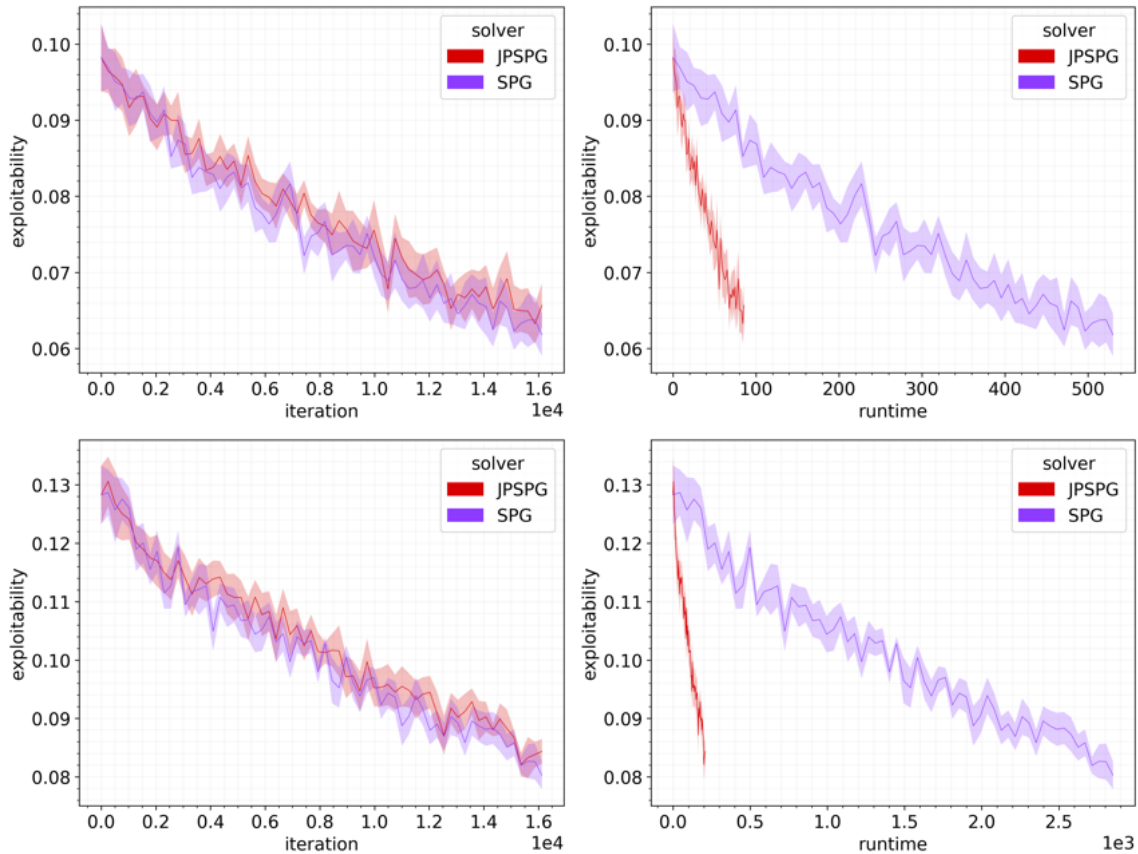


Figure 5.33: Unit-demand auction. Top: 10 players, 10 items. Bottom: 20 players, 20 items.

1859 5.3.2.2 Knapsack auction

1860 Another type of auction is a **knapsack auction**. We follow the description given in Aggarwal and
 1861 Hartline (2006). In a knapsack auction, an auctioneer auctions off space in a knapsack of known
 1862 capacity C . Each player seeks to place exactly one object in the knapsack. Player i values the
 1863 placement of its object in the knapsack at v_i . The valuations are private data of each respective
 1864 player. Each object takes up a certain amount of space in the knapsack. Player i 's object takes
 1865 space c_i , and these sizes are publicly known. Thus, the c_i s are public while the v_i s are private.
 1866 Each player submits a bid b_i . Among other works, knapsack auctions have been studied by Dütting,
 1867 Gkatzelis, and Roughgarden (2014) and Berg et al. (2010), who model the problem of bidding in
 1868 ad auctions as a penalized multiple choice knapsack problem. As Aggarwal and Hartline (2006)
 1869 note, the knapsack auction problem models several interesting applications. For example, consider
 1870 running a single auction to sell advertising space on a web page over the course of a day. Suppose
 1871 statistical information is available for each advertiser as to how many showings (i.e., impressions)
 1872 are necessary to result in a user click-through and as well as how many times the web page itself will

1873 be viewed in a day. The number of impressions necessary to generate a click-through corresponds to
1874 the c_i s and the number of total views corresponds to the capacity of the knapsack, C . Each utility
1875 function evaluation requires solving an optimization problem of the following form: maximize $\mathbf{x} \cdot \mathbf{b}$
1876 subject to $\mathbf{x} \cdot \mathbf{c} \leq C$ and $\mathbf{x} \in \{0, 1\}^n$. Here, n is the number of players, \mathbf{b} is the vector of stated
1877 values (bids) for each player, \mathbf{c} is the corresponding vector of sizes, \mathbf{x} is a binary vector indicating
1878 whether each player is included in the knapsack. Player i 's final utility is $(v_i - b_i)x_i$. That is, it is
1879 the difference between their private value and their submitted bid, assuming they are included in
1880 the knapsack, and zero otherwise. This problem can be solved using **integer linear programming**
1881 (ILP). For this, we use the `milp` function included in SciPy's library (Virtanen et al., 2020), which
1882 uses the HiGHS optimization solver (Huangfu and Hall, 2018; Hall et al., 2023). Solving an integer
1883 program can be expensive (integer linear programming is NP-hard), so reducing the number of
1884 utility function evaluations during learning should result in a significant speedup. In our experiment,
1885 we sample the v_i s and c_i s from the standard uniform distribution, and sample C from the standard
1886 uniform distribution on $[0, n]$. Experimental results on the knapsack auction are shown in Figure
1887 5.34. Our method requires fewer utility function evaluations per iteration and thus yields a dramatic
1888 improvement in training time for attaining the same exploitability.

1889 5.3.2.3 Sequential auction for identical items

1890 Consider a multi-item unit-demand auction in which identical items are sold *sequentially*, rather than
1891 simultaneously. We follow the description given in Krishna (2002, §15.1). In this auction, K identical
1892 items are sold to $N > K$ bidders using a sequence of first-price sealed-bid auctions. Specifically,
1893 on each of K rounds, one of the items is auctioned using the first-price format, and the price at
1894 which it is sold—the winning bid—is announced. We focus on the **single-unit demand** setting, in
1895 which each bidder has use for at most one unit. Thus a bidder leaves the game once it has won an
1896 item. Each bidder has a private value v_i that is sampled from the standard uniform distribution.
1897 On round k , a bidder's observation consists of its own private value as well as all the prices of the
1898 preceding $k - 1$ rounds, p_1, p_2, \dots, p_{k-1} . Results are shown in Figure 5.35. Our method yields a
1899 dramatic improvement in terms of the wall time required to reach a certain level of exploitability.

1900 5.3.2.4 Continuous-action Goofspiel

1901 Goofspiel, also known as *the game of pure strategy*, is a card game invented by mathematician Merrill
1902 Flood in the 1930s (Tucker, 1984). This game is played with a standard card deck. The cards of one
1903 suit are given to one player, the cards of a second suit are given to the other player, and the cards
1904 of a third suit are shuffled and placed face down in the middle. The cards are valued from low to
1905 high as 1 (Ace), 2, 3, \dots , 10, 11 (Jack), 12 (Queen), and 13 (King). A round consists of turning
1906 up the next card from the middle pile and letting the players simultaneously “bid” on this “prize”
1907 card. Players bid by choosing one of their own cards and revealing it at the same time as the other
1908 player. The player with the highest bid wins the value of the prize card. In a tie, the prize value
1909 is split between the players. All three cards are then discarded. The game ends after 13 rounds,
1910 and the winner is the player with the highest score. Because of its simple mechanics but complex
1911 strategy, Goofspiel is commonly used as an example in game theory and artificial intelligence, and
1912 has been studied extensively in the literature (Ross, 1971; Dror, 1989; Ferguson and Melolidakis,
1913 2001; Grimes and Dror, 2013; Rhoads and Bartholdi, 2012; Lanctot, Lisý, and Winands, 2014). We
1914 consider the following *continuous-action* variant of Goofspiel. Instead of receiving a deck of discrete
1915 bids consisting of all cards of one suit, each player has a *continuous* budget that they can spend to

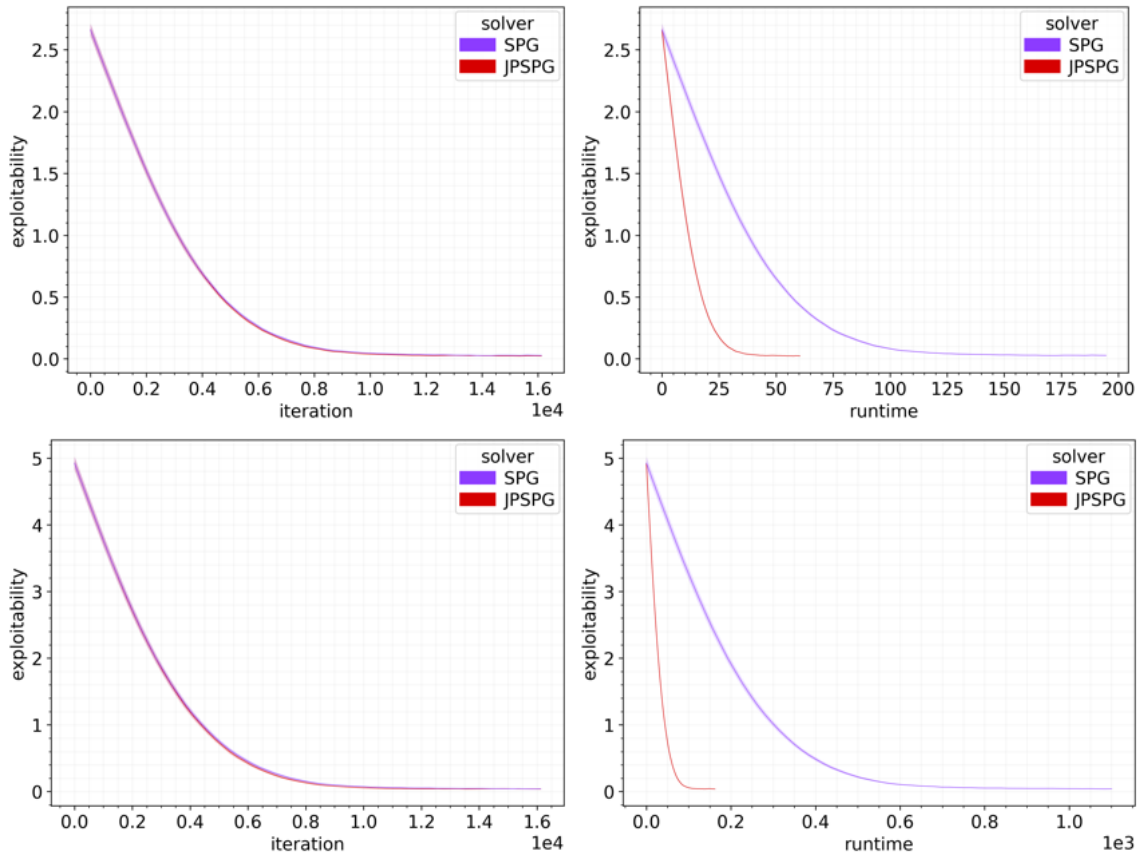


Figure 5.34: Knapsack auction. Top: 10 players. Bottom: 20 players.

1916 bid on the prize card in each round. In each round, their continuous bid is subtracted from their
 1917 budget. This can be thought of as a multi-round, multi-item, auction-like scenario with a budget
 1918 constraint for each bidder. To allow the players to randomize over their 1-dimensional actions (the
 1919 bids), we inject their strategy networks with 1-dimensional latent input noise in addition to the
 1920 observation, as described in Martin and Sandholm (2023). Figure 5.36 shows the exploitability over
 1921 the course of training on continuous-action Goofspiel. Our method yields a significant improvement
 1922 in the run time required to attain a certain level of exploitability.

1923 Our experimental results confirm our hypothesis, namely, that our approach yields a dramatic
 1924 improvement in the training time required to reach a certain level of exploitability.

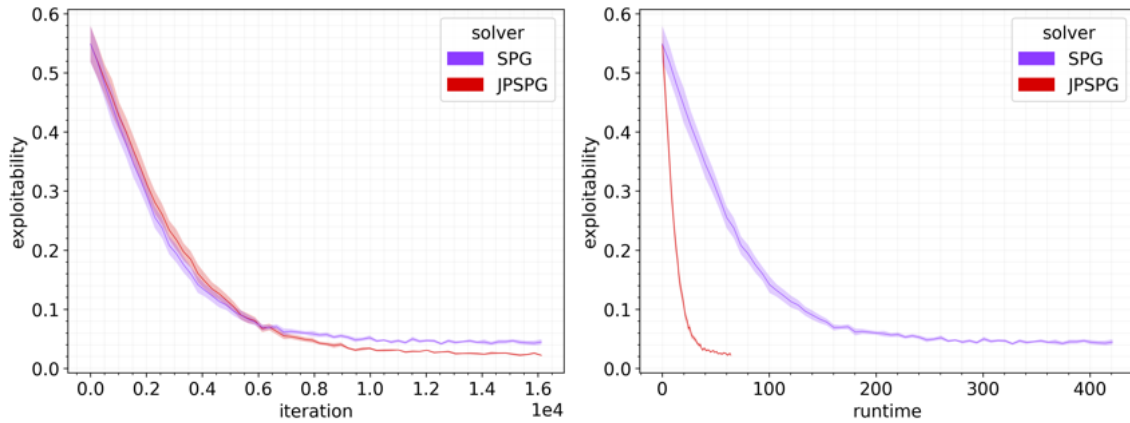


Figure 5.35: 20-player, 10-item sequential auction.

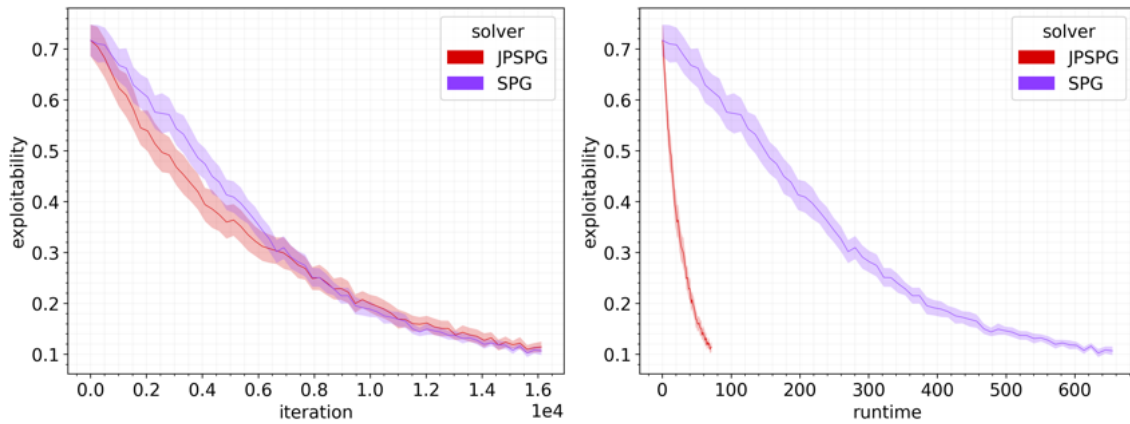


Figure 5.36: 20-player continuous Goofspiel.

1925 5.4 Solving infinite-player games with player-to-strategy net- 1926 works

1927 In Martin and Sandholm (2025e), we present a new approach to solving games with a countably or
1928 uncountably infinite number of players. Such games are often used to model multiagent systems
1929 with a large number of agents. The latter are frequently encountered in economics, financial markets,
1930 crowd dynamics, congestion analysis, epidemiology, and population ecology, among other fields. Our
1931 two primary contributions are as follows. First, we present a way to represent strategy profiles for
1932 an infinite number of players, which we name a Player-to-Strategy Network (P2SN). Such a network
1933 maps players to strategies, and exploits the generalization capabilities of neural networks to learn

1934 across an infinite number of inputs (players) simultaneously. Second, we present an algorithm for
 1935 training such a network to find approximate Nash equilibria, which we name Shared-Parameter
 1936 Simultaneous Gradient (SPSG). This algorithm generalizes simultaneous gradient ascent and its
 1937 variants, which are traditionally used for multiagent reinforcement learning. We test our approach
 1938 on infinite-player games and observe its convergence to approximate Nash equilibria. Our method
 1939 can handle games with infinitely many players, infinitely many actions (and mixed strategies over
 1940 them), infinitely many states, and discontinuous player-specific utility functions. To our knowledge,
 1941 this is the first paper to tackle general infinite-player games using neural networks.

1942 5.4.1 Method

1943 Our method consists mainly of two parts: Strategy profile representation and training.

1944 **P2SN.** Suppose we are interested in tackling a game with infinitely many players. This raises the
 1945 problem of how to represent a strategy profile that has an infinite number of players. We propose a
 1946 way to do this, which we call a **Player-to-Strategy Network** (P2SN). This is a neural network that
 1947 takes as input a *player* and outputs a *strategy* for that player. We exploit the strong generalization
 1948 capabilities of neural networks (as described in Section 1.1) to represent (and potentially learn, as
 1949 we will see shortly) across an infinite number of possible inputs, i.e., players, simultaneously.⁶ What
 1950 exactly the network receives as input depends on the game. The inputs include (1) the features
 1951 identifying the player, (2) any observations the player receives, and (3) random noise that allows the
 1952 network to randomize over actions (Martin and Sandholm, 2023). This is illustrated in Figure 5.37.

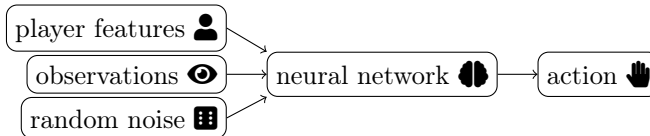


Figure 5.37: Structure of a P2SN.

1953 **Player representation.** A specific player is input into the P2SN as a set of *player features*.
 1954 The exact nature of these features depends on the game. For example, in a spatial game, each
 1955 player might be identified with a specific point in two-dimensional space. Thus a player can be
 1956 represented as an ordered pair of X and Y coordinates, i.e., a 2D vector. In a game with a continuum
 1957 of indistinguishable players, such as a continuum of traders in a stock market (Aumann, 1964), the
 1958 players can be represented simply as points on the unit interval (i.e., scalars), with no additional
 1959 distinguishing features. In other games, if there is some prior notion of *similarity* between players,
 1960 they can be represented as points in some embedding space, such that similar players get nearby
 1961 embeddings and dissimilar players get distant embeddings. This type of learning, which pulls together
 1962 similar objects and pushes apart dissimilar objects, is known as contrastive learning (Jaiswal et al.,
 1963 2020; Le-Khac, Healy, and Smeaton, 2020; Hu, Wang, et al., 2024).

⁶P2SN could be useful even in games with a finite number of players. For instance, the input could be a one-hot encoding of the player’s index, along with any other *a priori* features of that player. It allows for neural net parameters to be shared across players, while simultaneously taking into account each player’s individual and specific *a priori* features as input. This can potentially result in faster learning and better generalization. However, the present paper focuses on games with an infinite number of players.

1964 **Fourier features.** The space of player feature vectors is often low-dimensional. For example,
 1965 it is 2-dimensional in the case of a spatial game on the square. It is important for our P2SN to
 1966 be able to represent fine details on such a space. It has been found that standard feedforward
 1967 neural networks struggle to learn detailed, high-frequency features in low-dimensional input spaces.
 1968 Rahaman et al. (2019) showed that standard neural networks are biased towards low-frequency
 1969 functions, meaning that they cannot have local fluctuations without affecting their global behavior.
 1970 More precisely, using Fourier analysis, they showed that these networks prioritize learning the
 1971 low-frequency modes, a phenomenon they call *spectral bias*. Tancik et al. (2020) showed that passing
 1972 input points through a simple **Fourier feature mapping** enables a **multilayer perceptron** (MLP)
 1973 to learn high-frequency functions in low-dimensional problem domains, whereas a standard MLP
 1974 has impractically slow convergence to high-frequency signal components. To allow the P2SN to more
 1975 easily represent and learn fine details in low-dimensional spaces, which we found to help performance,
 1976 we pre-process the inputs $\mathbf{x} \in \mathbb{R}^d$ with the map $f(\mathbf{x}) = (\sin(\mathbf{a}), \cos(\mathbf{a}))$, where $\mathbf{a} = \mathbf{B}\mathbf{x} + \mathbf{b}$, $n \in \mathbb{N}$ is
 1977 the number of Fourier features, $\mathbf{B} \in \mathbb{R}^{d \times n}$ is a learned frequency matrix, and $\mathbf{b} \in \mathbb{R}^n$ is a learned
 1978 phase vector. Thus f maps the inputs to a high-dimensional space of Fourier features. We initialize \mathbf{B}
 1979 with independent samples from the normal distribution with standard deviation $\sigma = 100$. This large
 1980 σ is recommended by Tancik et al. (2020) and yields high spatial frequencies from the start, which
 1981 facilitates representation and learning. We initialize the entries of \mathbf{b} from the uniform distribution
 1982 on $[0, 2\pi)$.

1983 **Learning dynamics for shared parameters.** If parameters are shared between players (instead
 1984 of maintaining individual disjoint parameters for each player), the original expression for $v(s)_i$ no
 1985 longer makes sense. This is because it requires taking a derivative with respect to the parameters of
 1986 a specific player and no other, namely s_i , which is no longer possible.

1987 **Naive SG.** To get around this, we could, for each player, simply take derivatives with respect
 1988 to *all* parameters: $\tilde{v}(s)_i = \frac{du(s)_i}{ds}$. However, this does not work. For example, consider a two-player
 1989 zero-sum game where both players share parameters $s \in [0, 1]^2$, player 1’s strategy is s_1 , player 2’s
 1990 strategy is s_2 , and the utility for player 1 is $s_1 + s_2$. If we just ignore the disjointness requirement, the
 1991 parameters do not change at all, because they are “pulled” in exactly equal and opposite directions.
 1992 We demonstrate this in Section 5.4.2.1.

1993 **Shared-Parameter Simultaneous Gradient (SPSG).** To solve these problems, we re-express
 1994 the simultaneous gradient as a derivative with respect to *the entire strategy profile* as a unit:
 1995 $v(s) = \left[\frac{d}{dr} \int_{i \sim \mu} u(s[i \mapsto r_i])_i \right]_{r=s}$. Here, $[\cdot]_{r=s}$ means evaluating the expression inside the brackets,
 1996 treated as a function of r , at the argument value s . The derivative of the integral is a **functional**
 1997 **derivative**, i.e., the derivative of a functional (the integral) with respect to a function (in this
 1998 case, r). In contrast to the original expression for the simultaneous gradient, this expression *can* be
 1999 generalized to the shared-parameter case. Let $s \in \Theta \rightarrow \prod_{i \in P} S_i$ be a parameterization of a strategy
 2000 profile—such as a P2SN—with parameters in Θ . Define $v(\theta) = \left[\frac{d}{d\phi} \int_{i \sim \mu} u(s_\theta[i \mapsto (s_\phi)_i])_i \right]_{\phi=\theta}$. We
 2001 call this the **Shared-Parameter Simultaneous Gradient** (SPSG).

2002 **Implementation.** In practice, SPSG can be estimated as follows. On each iteration, do the
 2003 following. First, create a copy ϕ of the current parameters θ . Second, sample a player i from μ .

Algorithm 2 Training a P2SN with SPSG.

Input: Initial parameters $\theta \in \mathbb{R}^d$, learning rate $\eta \in \mathbb{R}$.

$f(\theta, \phi)_i = u(s_\theta[i \mapsto (s_\phi)_i])_i$

for $t \leftarrow 0, 1, 2, \dots$ **do**

$i \sim \mu$

\triangleright sample a player from the player distribution

$g \leftarrow \nabla_2 f(\theta, \theta)_i$

\triangleright get unbiased estimator of the SPSG

$\theta \leftarrow \theta + \eta g$

\triangleright update the shared parameters

2004 Third, create a hybrid strategy profile $s_\theta[i \mapsto (s_\phi)_i]$ that intertwines s_θ and s_ϕ , using s_ϕ 's output to
2005 define player i 's strategy and s_θ 's output to define any other player's strategy. Fourth, pass this
2006 strategy profile to the utility function $u(\cdot)_i$ and get player i 's utility. Fifth, take the gradient of this
2007 scalar with respect to ϕ .⁷ The result is an unbiased estimator of $v(\theta)$. The sampling of players yields
2008 an unbiased estimator of the entire integral, which in turn yields an unbiased estimator of the SPSG,
2009 since the integral and derivative commute.⁸ All we need for this procedure is access to $u(s)_i$, or an
2010 unbiased estimator thereof. This shown in Algorithm 2.

2011 **Estimating integrals.** In cases where $u(s)_i$ is an integral over players of some other function,
2012 which is sometimes the case in infinite-player games, we can use an unbiased estimator for it. Such
2013 an estimator often exists for infinite-player games in practice. For example, many infinite-player
2014 games can be expressed as a series of the following general form.

$$u(s)_i = f_i(s_i) + \int_{j \sim \nu_i} \left(g_{ij}(s_i, s_j) + \int_{k \sim \xi_{ij}} h_{ijk}(s_i, s_j, s_k) + \dots \right) \quad (5.27)$$

2015 where f, g , etc. are arbitrary functions and ν_i, ξ_{ij} , etc. are arbitrary measures. That is, player i 's
2016 utility is an integral over pairwise interactions with other players, 3-way interactions with other
2017 pairs of players, etc., up to some order. In that case, an unbiased estimator of player i 's utility can
2018 be obtained by first sampling j given i , then sampling k given i and j , and so on, for however many
2019 terms are necessary. If necessary, we can use techniques like **Markov chain Monte Carlo** (MCMC)
2020 (Metropolis et al., 1953; Hastings, 1970) to obtain these samples. As we will see in the experiments
2021 section, many games of interest require only the first integral, i.e., involve only the aggregation of
2022 pairwise interactions between players (and their corresponding strategies). We emphasize that our
2023 method makes no assumption of symmetry or identicality across players. Each player can have its
2024 own distinct strategy space and utility function. For example, in spatial games (where players occupy
2025 different points in space), the players are clearly not identical: some are close to the boundary (where
2026 boundary effects come into play), and others are not. Furthermore, the computational complexity of
2027 our method, when applied to the case where P is finite, is close to that of SGA, differing only in
2028 that a *single* player is *stochastically* sampled at each iteration (and used to estimate the SPSG).
2029 P2SN and SPSG leverage the generalization capabilities of neural networks to learn across an infinite
2030 number of players simultaneously.

⁷Not θ . Though θ and ϕ have equal *values*, they are different *variables*. Consider $f(x, y) = xy^2$. At $x = y = 1$, x and y have the same value. However, the gradient with respect to x is 1, while the gradient with respect to y is 2.

⁸Here, we assume that the conditions of the Leibniz integral rule, which allow differentiation and integration to be exchanged, hold, as described in Talvila (2001).

2031 5.4.2 Experiments

2032 We evaluated the proposed techniques using computational experiments. For each experiment, we
2033 ran 16 trials. In each plot of exploitability, solid lines show the mean across trials, and bands
2034 show a confidence level of 0.95 for the mean, which is computed using bootstrapping (Efron, 1979),
2035 specifically the **bias-corrected and accelerated** (BCa) method (Efron, 1987). The legend entries
2036 are sorted by final values. In the neural net, on each iteration, we use a batch size of 64, averaging
2037 gradients across that batch. For the P2SN, we use a multilayer perceptron with various depths
2038 (number of hidden layers) and widths (number of Fourier features, as well as size of each hidden
2039 layer). If not stated explicitly, we use a default depth of 1 and width of 64. We use the ReLU
2040 activation function and He weight initialization (He et al., 2015). We use a vanilla SGD optimizer
2041 with learning rate as stated in the legends (labeled “lr”).

2042 **Exploitability estimation.** To evaluate performance, we estimate the exploitability. To do this,
2043 we discretize the space of players into N_{players} points⁹, estimate the exploitability of each individual
2044 player, and average. To estimate an individual player’s exploitability, we estimate the quantities
2045 described in the formulation of Bayesian games at the end of Section 3.3. Specifically, we sample
2046 $N_{\text{observations}}$ observations from that player’s observation distribution to average over, discretize the
2047 action space into N_{actions} points to maximize over, and sample N_{samples} utilities conditioned on the
2048 observation to average over.¹⁰ We use $N = 256$ for each of these.

2049 For brevity, we include only some experiments here. The rest can be found in Martin and
2050 Sandholm (2025e).

2051 5.4.2.1 Sum game

2052 To illustrate why Naive SG fails, and why we need SPSG instead, we use a very simple game. This is
2053 a two-player zero-sum game in which both players simultaneously choose values on the unit interval,
2054 and the utility for the first player is the sum of the values. Formally, $P = \{1, 2\}$, $S(p) = [0, 1]$,
2055 $u(s, 1) = -u(s, 2) = s(1) + s(2)$. We run training with the two different dynamics. Exploitabilities
2056 are shown in Figure 5.38. As expected, the naive dynamics fail to converge, because the parameters
2057 are pulled in equal and opposite directions (with some noise from which player is sampled in each
2058 iteration). On the other hand, SPSG does converge.

2059 5.4.2.2 Anti-coordination game

2060 Let $d \in \mathbb{N}$ be a spatial dimension. Let $r \in \mathbb{R}$ be an interaction radius. Let $k \in \mathbb{N}$ be a number of
2061 resources. Let $b \in \mathbb{R}^d \rightarrow \mathbb{R}^k$ be a bias field. The player set is $P = [0, 1]^d$. The strategy set function is
2062 $S(x) = \Delta[k]$. The utility function is $u(s, x) = s(x) \cdot (b(x) - \bar{s}(x))$, where $\bar{s}(x) = \int_{y \in B_r(x) \cap P} s(y)$ is the
2063 aggregate strategy around x . This game models a situation where there are k congestible resources
2064 and a continuum of players arranged in space. Each player derives a certain value from using each
2065 resource, but their ability to use the resource diminishes with the amount of surrounding players who

⁹We discretize a 1D interval in the obvious way, i.e., by placing N equally-spaced points across the interval. We discretize a higher-dimensional cube by taking the first N points of the Roberts sequence (Roberts, 2018), a low-discrepancy (a.k.a. quasi-random) sequence (Kuipers and Niederreiter, 1974; Niederreiter, 1992) with good approximation properties. We discretize a simplex using a transformation of the low-discrepancy sequence for the cube (Pillards and Cools, 2005).

¹⁰More details, including how to sample states conditioned on observations, can be found in Bichler, Fichtl, Heidekrüger, et al. (2021).

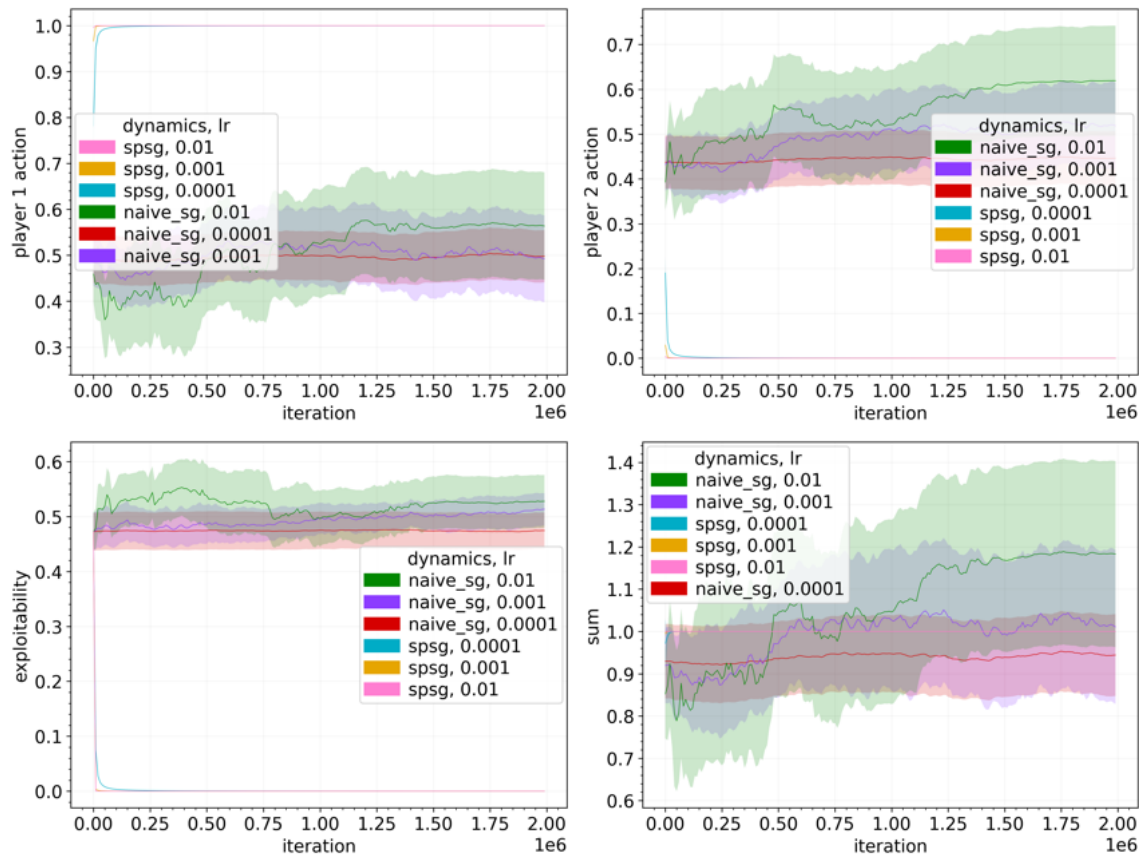


Figure 5.38: Sum game.

2066 *also* choose that resource. Thus each player wants to satisfy a particular endogenous *bias*, but also
 2067 wants to exogenously *anti-conform* or anti-coordinate with its neighbors. The game resembles the
 2068 Lenz–Ising model (Lenz, 1920; Ising, 1925) of statistical mechanics, which models ferromagnetism
 2069 and the magnetic dipole moments of atomic spins. In particular, each player’s strategy corresponds
 2070 to an **expected spin vector**, and the integral of the utility function over all players (i.e., the
 2071 utilitarian social welfare) yields the negative Hamiltonian of that model. This so-called **Ising game**
 2072 has been studied by Galam and Walliser (2010), Xin, Yang, and Huang (2017), Leonidov, Savvateev,
 2073 and Semenov (2020), Leonidov, Savvateev, and Semenov (2024), and Feldman, Kim, and Palmer
 2074 (2024), among others. In our experiments, we use $r = 0.1$, $k = 3$, and $b(x) = [x \in B_{1/4}(\mathbf{1}/2)]\mathbf{e}_0$.
 2075 The latter incentivizes players near the center of the domain to choose 0. We do this to break
 2076 spatial symmetry and eliminate a trivial everywhere-constant equilibrium where players mix equally
 2077 between all resources. Exploitabilities are shown in Figure 5.39 and 5.40. In both cases, our method
 2078 decreases exploitability across iterations. Learned strategy profiles are shown in Figure 5.39 and 5.40.
 2079 As expected, players inside the ball at the center of the domain choose 0. Beyond this ball, players

2080 alternate between the 3 choices (except immediately around the ball, where the players always skip
 2081 choice 0).

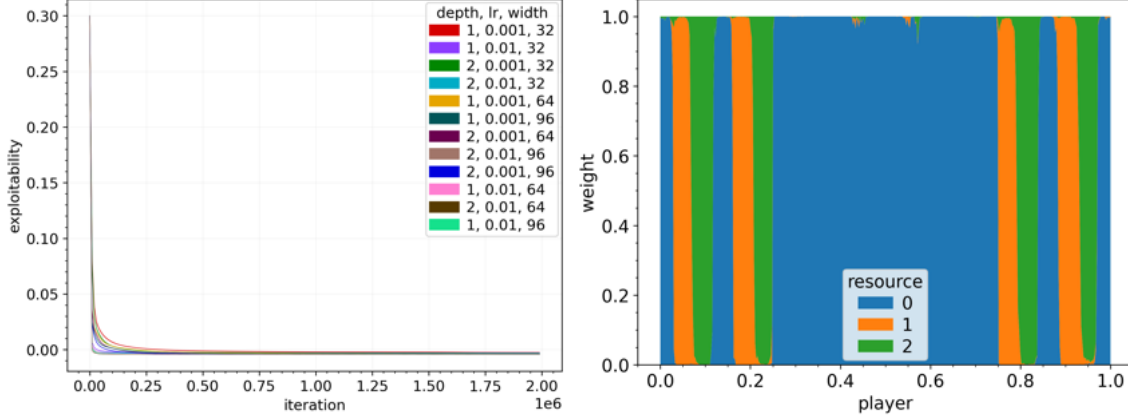


Figure 5.39: 1D anti-coordination game. Left: Exploitabilities. Right: A learned strategy profile. This is a stacked bar plot: each vertical slice shows a point on the simplex.

2082 5.4.2.3 Discontinuous game

2083 Let $\varepsilon, \delta > 0$. Let $P = [0, 1]$, $S(p) = [0, 1]$, and $u(s, p) = \int_{q \in B_\varepsilon(p) \cap P} [|s(p) - s(q)| < \delta]$ for $p \in$
 2084 $[0.1, 0.4] \cup [0.6, 0.9]$, $s(p)$ for $p \in [0, 0.1] \cup (0.9, 1]$, and $-s(p)$ for $p \in (0.4, 0.6)$. Each player in the
 2085 first region has an incentive to choose an action within δ of its neighbors', making this a type of
 2086 conformity game. The integrand is discontinuous. Consequently, taking the gradient of the integrand
 2087 yields an incorrect, biased estimator of the gradient of the integral. We overcome this obstacle by
 2088 replacing the gradient used inside the SPSG with a pseudo-gradient, as described in Section 2.8.
 2089 In our experiments, we use $\varepsilon = \delta = 0.1$. Exploitabilities are shown in Figure 5.41. In the legend,
 2090 "smooth_scale" denotes the smoothing scale σ used for the pseudo-gradient, with $\sigma = 0$ reserved for
 2091 no smoothing. As expected, smoothing is needed to learn a strategy profile with low exploitability.

2092 5.4.2.4 Circle game

2093 Let S^1 be the unit circle. Let $r \in (0, \frac{1}{2})$ be an interaction radius. Let $P = S^1$, $S(p) = S^1$, and
 2094 $u(s, p) = 2 \int_{q \in [p-r, p]} d(s(p), s(q)) - \int_{q \in [p, p+r]} d(s(p), s(q))$. Here, $d(x, y) = \min\{|x - y|, 1 - |x - y|\}$
 2095 denotes the distance on the unit circle (with wraparound). The utility contains a repulsion term and
 2096 an attraction term. Each player wants to be similar to the players ahead of it, but more dissimilar
 2097 to the players behind it. Around a circle, that creates a chasing effect: everyone wants to pull away
 2098 from those just behind, which in turn makes those behind want to pull away from someone else, and
 2099 so on. With stronger repulsion than attraction, this chase can never close into a consistent profile.
 2100 Thus this game has no pure-strategy NE.¹¹ Therefore, players need to be able to randomize over the
 2101 continuous action space. To do this—i.e., model mixed strategies over the infinite action space—we

¹¹It has an MSNE where each player mixes uniformly on S^1 .

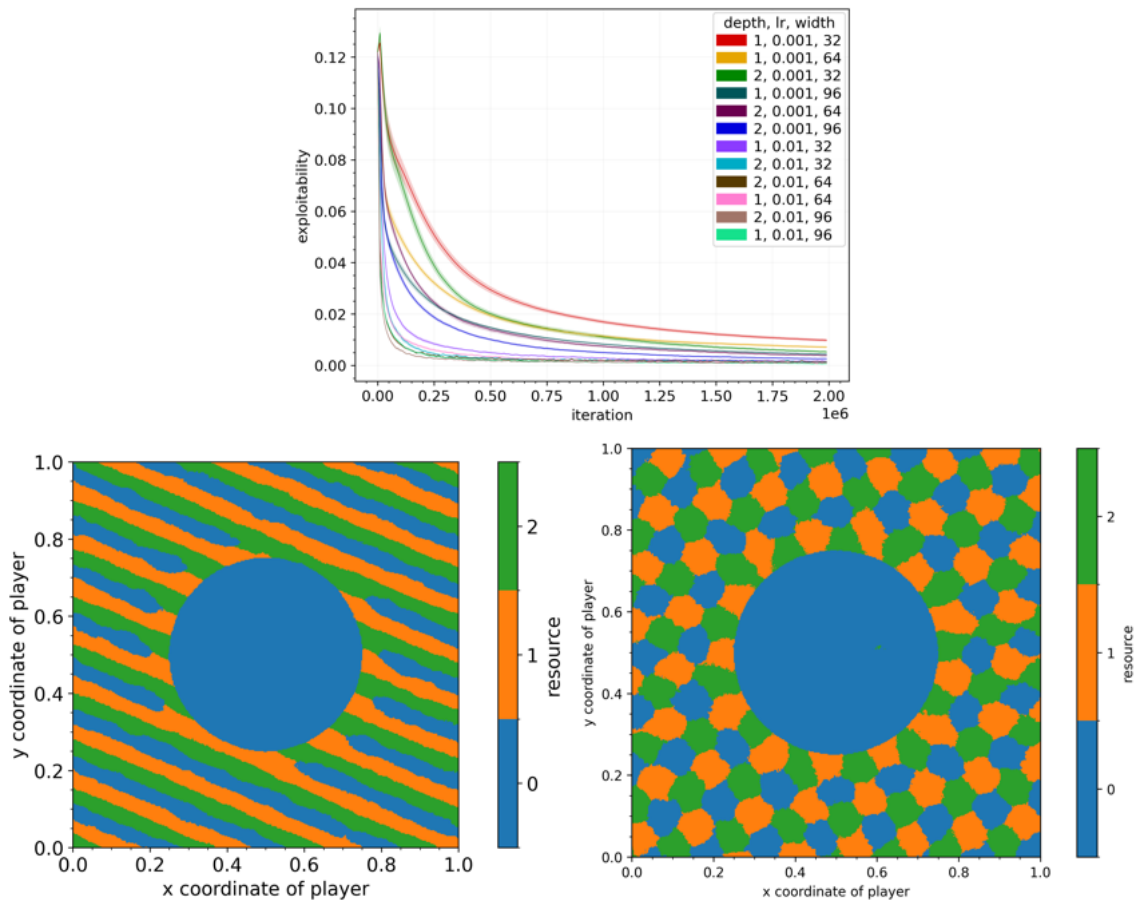


Figure 5.40: 2D anti-coordination game. Top: Exploitabilities. Bottom: Two learned strategy profiles. This is a categorical plot: each point shows the player’s top choice.

2102 use the randomized policy network approach described by Martin and Sandholm (2023), which
 2103 injects latent noise into the input of the network to induce a randomized output. In our experiments,
 2104 we use $r = 0.1$. Exploitabilities are shown in Figure 5.42. In the legend, “noise dim” denotes the
 2105 dimensionality of the Gaussian noise passed as additional input to the network. Learned strategy
 2106 profiles are shown in Figure 5.43. As expected, noise is needed to learn a strategy profile that can
 2107 randomize.

2108 5.4.2.5 Cournot game

2109 Let $d \in \mathbb{N}$ be a spatial dimension. Let $r \in \mathbb{R}$ be an interaction radius. The player set is $P = [0, 1]^d$.
 2110 The strategy set function is $S(x) = [0, 1]$. The utility function is $u(s, x) = s(x)P(\bar{s}(x)) - C(x, s(x))$,
 2111 where $\bar{s}(x) = \int_{y \in B_r(x) \cap P} s(y)$ is the aggregate strategy around x . This game models a spatial

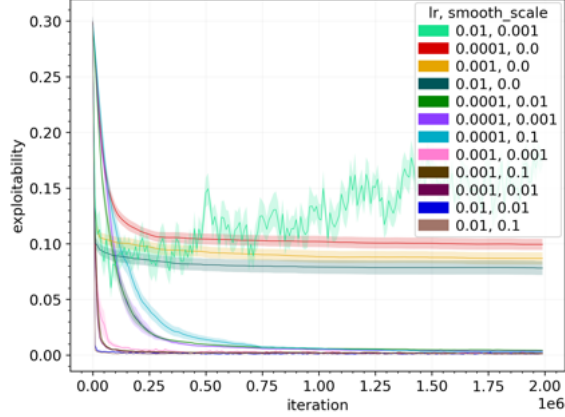


Figure 5.41: Discontinuous game.

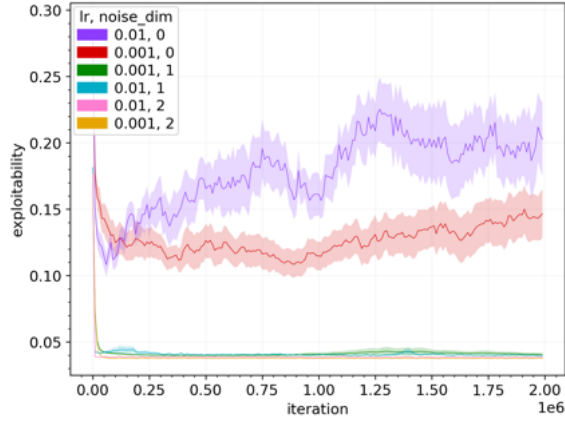


Figure 5.42: Circle game.

2112 Cournot competition. Cournot competition is a classic economic model dating back to Cournot
 2113 (1838) and Cournot (1863) in which firms compete in a market for a homogeneous good by choosing
 2114 a quantity of output to produce independently and simultaneously. Versions with a continuum of
 2115 players have been studied extensively (e.g., Novshek (1985) and Chan and Sircar (2015)). Each
 2116 player represents a *firm* choosing what quantity, if any, to produce of a good. A firm's profit is
 2117 revenue – cost = $qP(\bar{q}) - C(q)$, where q is its output, \bar{q} is the aggregate output of its neighborhood,
 2118 P is the *inverse demand function* (which yields the market price associated with an aggregate
 2119 output), and C is the firm's *production cost function*. For simplicity, we use a linear inverse demand
 2120 function $P(\bar{q}) = a - b\bar{q}$ and linear production cost function $C(q) = cq$. Thus, c represents a *marginal*
 2121 *cost of production*. In our experiments, we use $r = 0.1$, $a = b = 1$, and $c = 0.5$. Exploitabilities are
 2122 shown in Figure 5.44 and Figure 5.45. In both cases, our method decreases exploitability across
 2123 iterations. Learned strategy profiles are shown in Figure 5.44 and Figure 5.45. At the edges of the

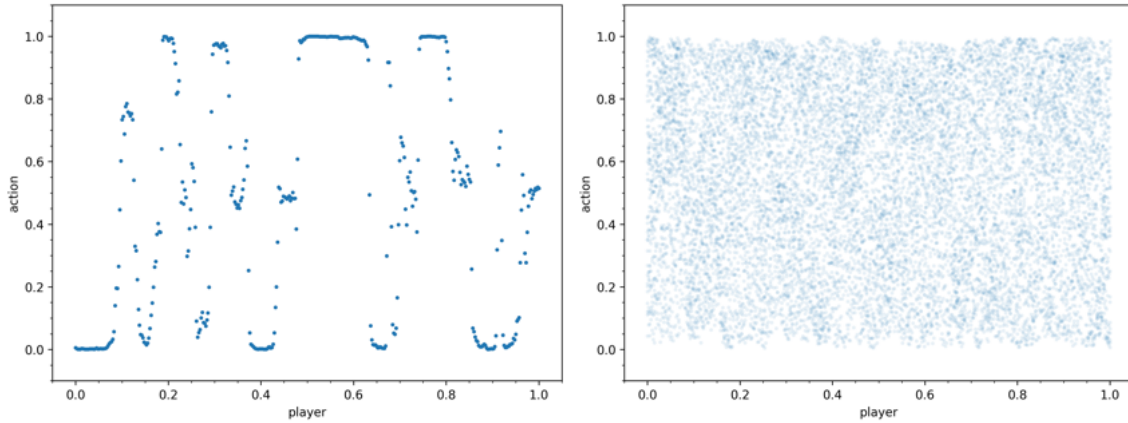


Figure 5.43: Learned strategy profiles for the circle game. Left: Noise dimension 0 (no noise). Right: Noise dimension 1.

2124 domain, players always choose to produce. This is because they are surrounded by fewer players in
 2125 total, and thus face less competition. Meanwhile, players in the interior of the domain alternate
 2126 between producing or not in narrow horizontal stripes.

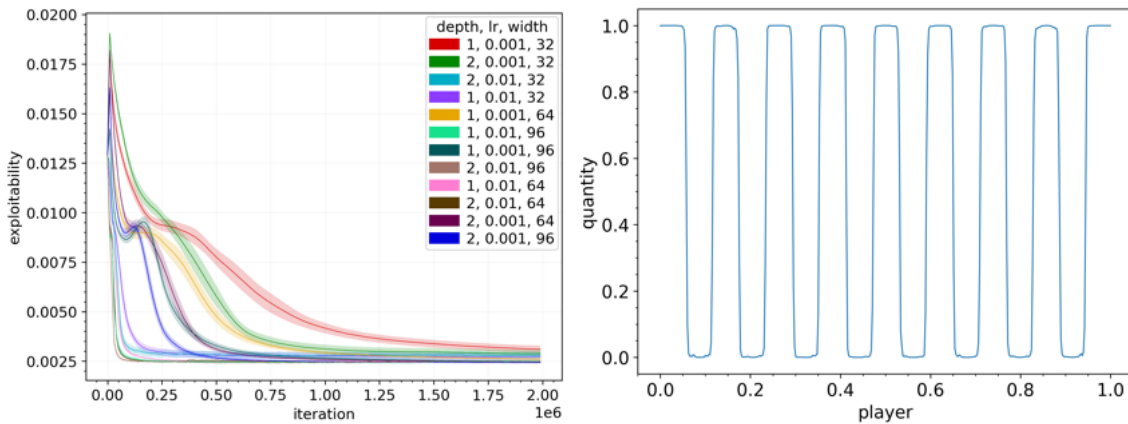


Figure 5.44: 1D Cournot game. Left: Exploitabilities. Right: A learned strategy profile.

2127 5.4.2.6 Bayesian Cournot game

2128 We now consider a variant of the Cournot game with **incomplete information**. Specifically, we
 2129 give each player only limited information about the marginal cost c . For this game, we let the state
 2130 be the marginal cost c , the state space be the unit interval $[0, 1]$, the initial state distribution be

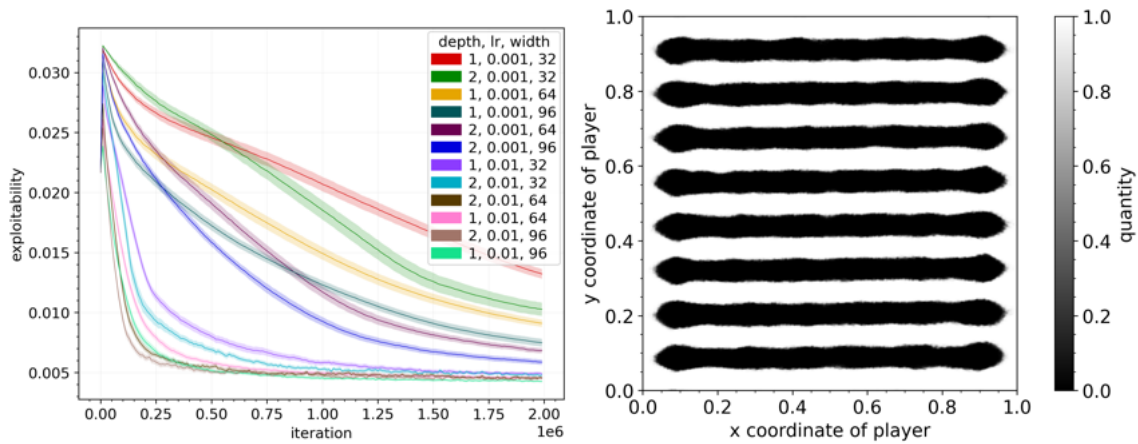


Figure 5.45: 2D Cournot game. Left: Exploitabilities. Right: A learned strategy profile.

2131 the standard uniform distribution, the observation space be \mathbb{R} , and the observation function for
 2132 player i yield $\max\{c, \|i\|_2\}$ deterministically. The actions and utility function are as before. This
 2133 game gives players closer to the origin more information. Exploitabilities are shown in Figure 5.46
 2134 for $d = 2$. Our method decreases exploitability across iterations. A learned strategy profile is shown
 2135 in Figure 5.46. As expected, players near the origin can make a better decision, because they have
 2136 more information about the true state of the world. For example, players near the origin decrease
 2137 production when the marginal cost is high.

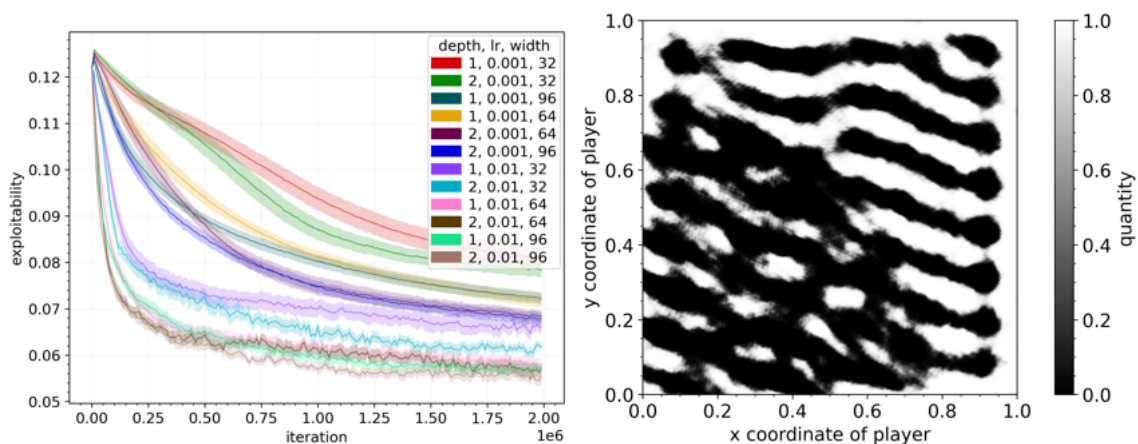


Figure 5.46: 2D Bayesian Cournot game. Left: Exploitabilities. Right: A learned strategy profile, under a state (marginal cost) of 0.8.

2138 **5.4.2.7 Quadratic-cost Cournot game**

2139 We replace the linear cost function $C(q) = cq$, with a *quadratic* cost function $C(q) = cq^2$. We keep
 2140 the same $c = 0.5$. Results are shown in Figure 5.47 and 5.48. Unlike the original Cournot game,
 2141 this time, most players choose outputs strictly between 0 and 1, with players in the interior of the
 2142 domain choosing values near 0.5. Players near the edges of the domain are, as before, more likely to
 2143 choose high outputs (of around 0.7 to 0.8), because they face less competition.

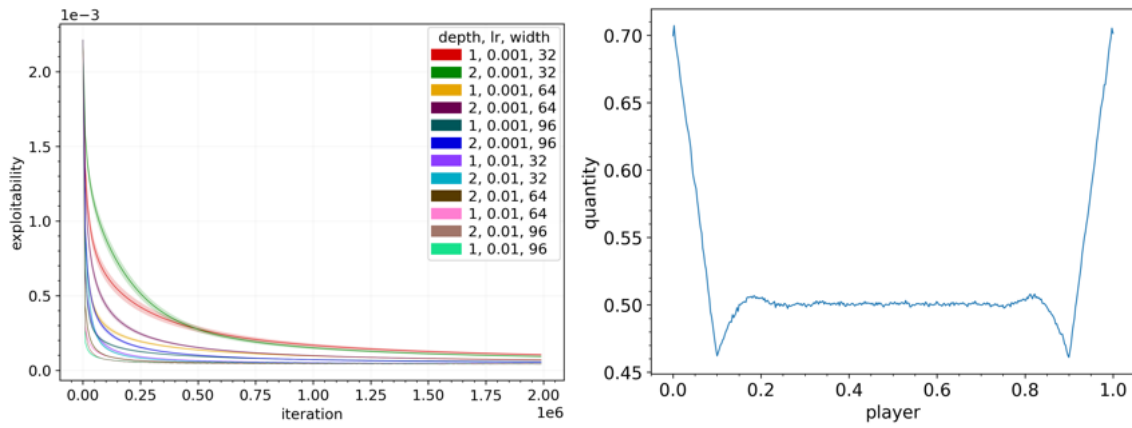


Figure 5.47: 1D quadratic-cost Cournot game. Left: Exploitabilities. Right: A learned strategy profile.

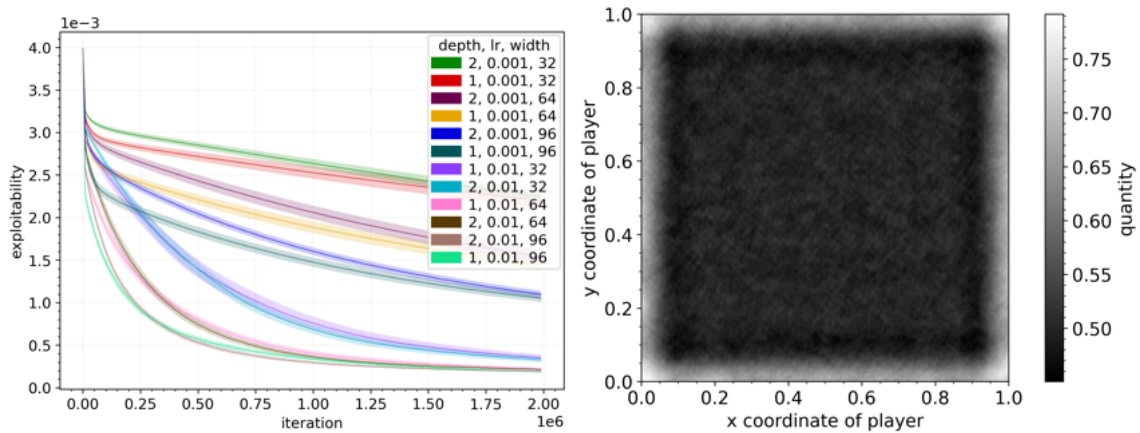


Figure 5.48: 2D quadratic-cost Cournot game. Left: Exploitabilities. Right: A learned strategy profile.

2144 **5.4.2.8 Heterogeneous-cost Cournot game**

2145 We now consider another variant of the spatial Cournot game. This time, we let the marginal cost c
 2146 vary across players. Specifically, we let $c(x, y) = \sin(\pi x)^2 \sin(\pi y)^2 + \frac{1}{2} \sin(3\pi x)^2 \sin(3\pi y)^2$. This is
 2147 shown in Figure 5.49. Exploitabilities are shown in Figure 5.50. A learned strategy profile is shown
 2148 in Figure 5.50. The strategy profile is heterogeneous in a way that reflects the underlying marginal
 2149 cost field.

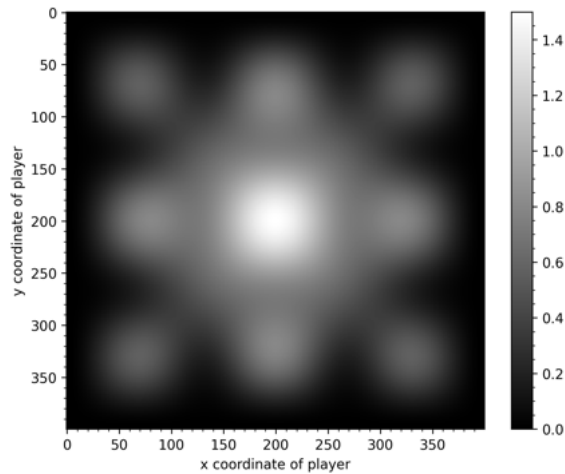


Figure 5.49: Marginal costs for heterogeneous-cost Cournot game.

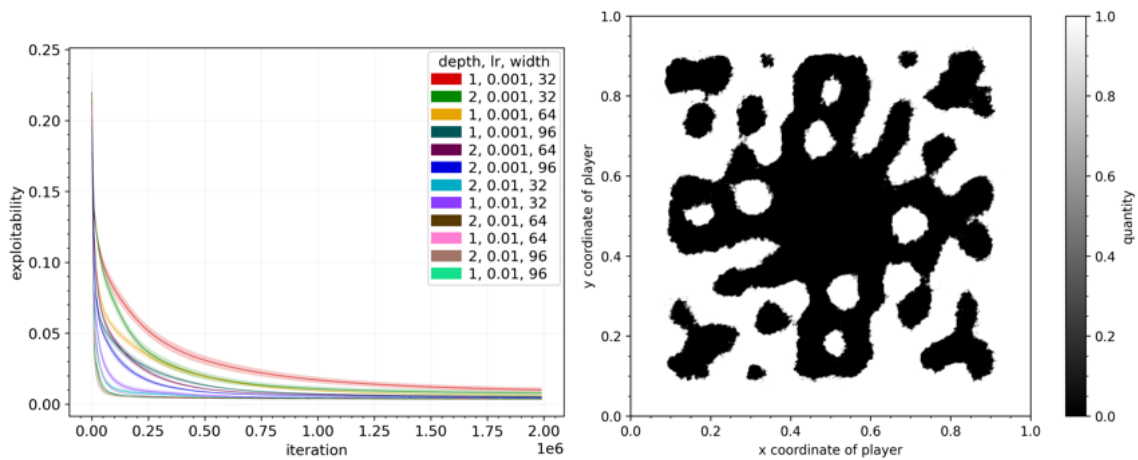


Figure 5.50: Heterogeneous-cost Cournot game. Left: Exploitabilities. Right: A learned strategy profile.

2150 **5.4.2.9 Conformity game**

2151 Let $d \in \mathbb{N}$ be a spatial dimension. Let $r \in \mathbb{R}$ be an interaction radius. The player set is $P = [0, 1]^d$.
 2152 The strategy set function is $S(x) = [0, 1]$. The utility function is

$$u(s, x) = - \begin{cases} (s(x) - 0)^2 & \max(x) < 0.3 \\ (s(x) - 1)^2 & \min(x) > 0.7 \\ (s(x) - \bar{s}(x))^2 & \text{otherwise} \end{cases} \quad (5.28)$$

2153 where $\bar{s}(x) = \int_{y \in B_r(x) \cap P} s(y)$ is the aggregate strategy around x . This models a situation where some
 2154 players (leaders) want to conform to a fixed target, while other players (followers) want to conform
 2155 to their neighbors. We let $r = 0.1$ in our experiment. Experimental results are shown in Figure 5.51
 2156 and Figure 5.52. As expected, leaders choose their target values, while followers interpolate between
 2157 their neighbors' values.

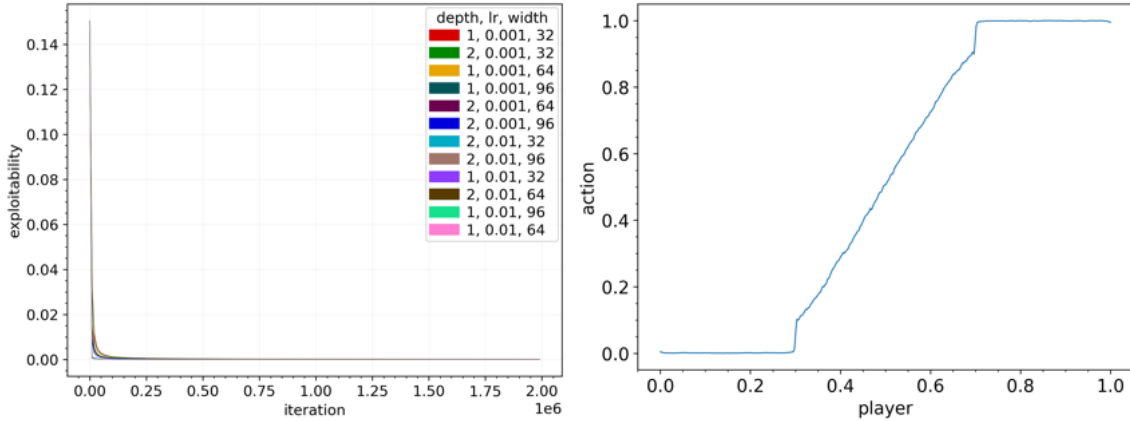


Figure 5.51: 1D conformity game. Left: Exploitabilities. Right: A learned strategy profile.

2158 **5.4.2.10 Permutation game**

2159 The **Plackett–Luce (PL) model** (Luce, 1959; Plackett, 1975) is a probability model for *rankings*,
 2160 that is, permutations of a set of items. It is based on Luce’s axiom of choice (Luce, 1959; Luce,
 2161 1977), which states that the probability of choosing one item over another does not depend on the
 2162 set of items from which the choice is made. Under this model, each item i has a score $x_i \in \mathbb{R}$, and
 2163 the probability of choosing item i from a set of items S is $\frac{\exp x_i}{\sum_{j \in S} \exp x_j} = \text{softmax}(\mathbf{x})_j$. A ranking
 2164 is modeled as a sequence of choices. At each point in the sequence, an item is chosen from the
 2165 remaining set of items according to the aforementioned probability. There is a fast, efficient way to
 2166 sample from the PL model, described in Yellott (1977), Grover et al. (2019, Proposition 5), and
 2167 Gadetsky et al. (2020, Lemma 1). Let $n \in \mathbb{N}$ be the number of items. For each item $i \in [n]$, let
 2168 $x_i \in \mathbb{R}$ be its score and $g_i \sim \text{Gumbel}$ be noise sampled from the standard Gumbel distribution.
 2169 Let the output permutation be $\pi = \text{argsort}(\mathbf{x} + \mathbf{g})$, where argsort yields the indices that sort its

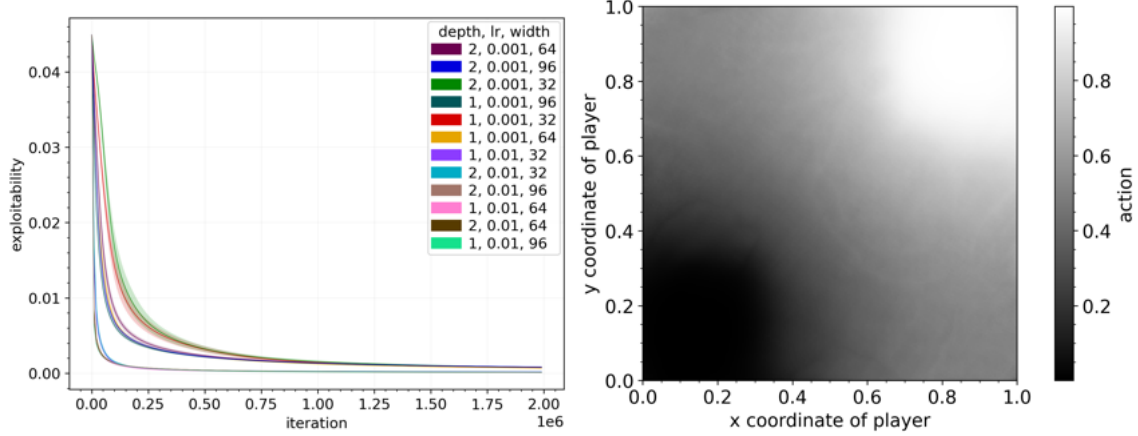


Figure 5.52: 2D conformity game. Left: Exploitabilities. Right: A learned strategy profile.

2170 input in descending order. This is sometimes called the *Gumbel argsort trick*. We denote the set of
 2171 permutations of $[n]$ by $n!$, so $\pi \in n!$ is an element of this set.

2172 Let $d, n \in \mathbb{N}$. Let $P = [0, 1]^d$ be the set of players. For each player $p \in P$, let $S(p) = \Delta[n]$ be the
 2173 set of strategies. Given a strategy profile s , for each player $p \in P$, let $\pi(p) \sim \text{PL}(\{\log s(p)_i\}_{i \in [n]})$
 2174 be a permutation sampled from its strategy according to the PL model described above (with the
 2175 samples being independent across players). In other words, each player independently chooses a
 2176 permutation of $[n]$. Therefore, this is a game with a combinatorial action space. Let the utility
 2177 function, expressed in terms of the sampled permutations, be

$$u(\pi, p) = \begin{cases} [\pi(p) = \text{id}_{[n]}] & \|p - \frac{1}{2}\|_2 < \delta \\ \int_{q \in B_r(p)} f(\pi(p), \pi(q)) & \text{otherwise} \end{cases} \quad (5.29)$$

2178 Here, $\delta, r \in \mathbb{R}_{\geq 0}$ are distances, $\text{id}_{[n]} = \{i\}_{i \in [n]}$ is the identity permutation, and $f : n! \times n! \rightarrow \mathbb{R}$ is a
 2179 pairwise interaction function. This incentivizes players inside a δ -ball at the center of the domain to
 2180 choose the identity permutation. Let the interaction function be as follows:

$$f(\pi, \sigma) = \sum_{i \in [n]} [\pi_i > \sigma_i] - [\pi_i < \sigma_i] \quad (5.30)$$

2181 That is, a player's payoff is the number of positions where it beats the opponent, minus the number
 2182 of positions where it loses to the opponent. This induces a non-transitive dominance relation between
 2183 permutations. Every permutation of length ≥ 3 is strictly dominated by another distinct permutation.

2184 We run experiments with $r = 0.1$, $\delta = 0.2$, and $d = 2$. Since this game has a discontinuous utility
 2185 function, we use pseudo-gradients. Exploitabilities are shown in Figure 5.53. A learned strategy
 2186 profile is shown in Figure 5.54.

2187 In Section C.2, we compare our method to the naive approach of discretizing the player space,
 2188 showing superior performance.

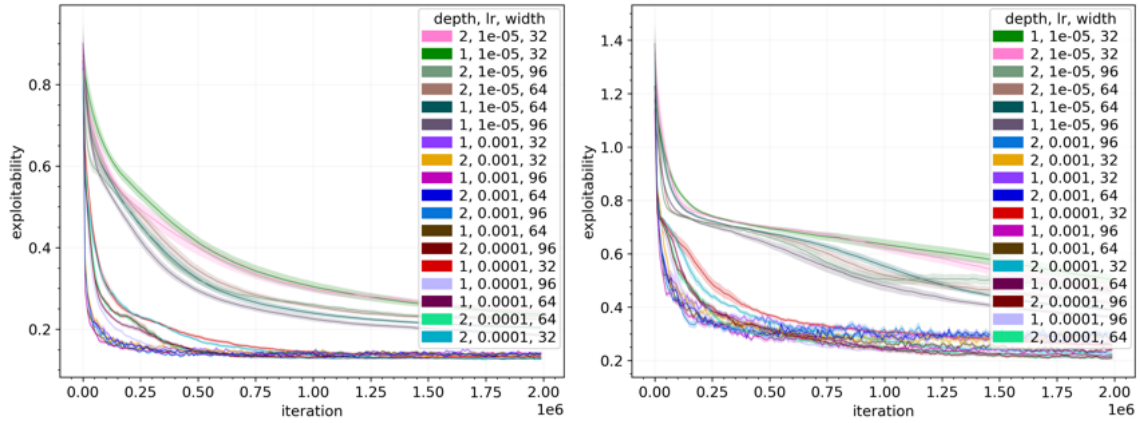


Figure 5.53: Permutation game with $n = 3$ and $n = 4$.

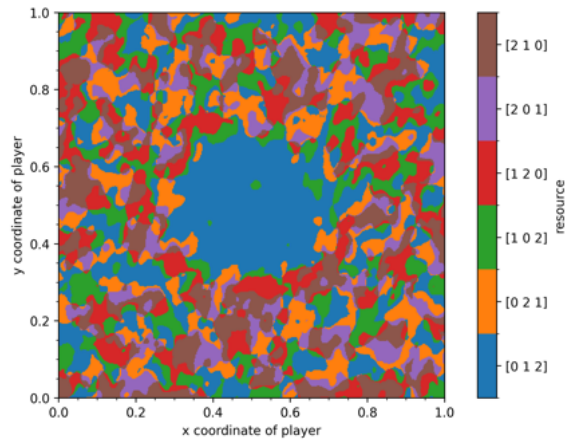


Figure 5.54: Strategy profile for permutation game with $n = 3$. This is a categorical plot. Each point shows the player's top choice.

2189 Chapter 6

2190 Proposed work

2191 In this section, we describe the next steps for our project as well as proposed future work. We aim
2192 to combine all of the techniques we have presented so far, to create a unified method that can
2193 solve (infinite-state, infinite-action, infinite-player, etc.) games in full generality. These results are
2194 obtainable within a reasonable amount of time. We want to emphasize that the methodological core
2195 is largely complete. The remaining work is primarily empirical and expository.

2196 6.1 Methods

2197 Our first step is to simply **integrate and unify** all of the techniques we have presented so far. We
2198 do this as follows:

- 2199 • Use a P2SN to represent strategy profiles.
- 2200 • Use SPSG for the simultaneous gradient analogue in this shared-parameter setting.
- 2201 • Use the infinite-player analogue of JPSPG for efficient estimation of this SPSG.
- 2202 • Use ApproxED-BRF for the learning dynamics.

2203 Further details and extensions are given below. We will test these on the proposed benchmarks,
2204 which are described in Section 6.2.

2205 6.1.1 Extension to multi-step games

2206 Our experiments with P2SN in the original paper (Section 5.4.1) were limited to single-step normal-
2207 form games and Bayesian games. However, we can straightforwardly generalize our method to apply
2208 to multi-step games. To do this, we can simply apply pseudo-gradients (Section 2.8) to the stochastic
2209 function that maps the parameters of the P2SN to sampled episode returns for a given player.

2210 In the multi-step case, we can no longer estimate best response values as we did for Bayesian
2211 games, that is, via the discretization we described at the end of Section 3.3. However, we can estimate
2212 best response values by running single-agent RL for the given player, treating the other players as a
2213 fixed part of the environment. To do this, we can use pseudo-gradients like we do in the multi-player
2214 case, but with just 1 player. This is essentially the technique used by Salimans et al. (2017). In
2215 particular, we block propagation of gradients through P2SN queries for other players. This is known

2216 as a “stop gradient” operation in automatic differentiation. It is the same operation we use to block
2217 gradients for other players when estimating the SPSPG.

2218 6.1.2 Extension to infinite-step games

2219 Our multi-step extension (Section 6.1.1) applies to games with a finite number of discrete timesteps.
2220 We propose to extend the framework further to the **infinite-step** regime—specifically, continuous-
2221 time **differential games** (Isaacs, 1965; Friedman, 1971a).

2222 In differential games, each player’s action at any point in time is a *rate* (velocity, thrust, bid
2223 intensity). For example, in a game where players control positions, their outputs could be a velocity
2224 (derivative of position with respect to time) instead of a displacement (change in position). Some of
2225 our proposed benchmarks, such as Orbit, Traffic, and Foraging, are naturally continuous-time, and
2226 treating them as genuine differential games (rather than coarse-grained discrete approximations)
2227 requires this extension.

2228 Differential games combine optimal control theory with game theory to model continuous-time
2229 conflicts. In these dynamic models, the state of the system evolves over time according to a set of
2230 (potentially stochastic) ODEs. The field was pioneered by Rufus Isaacs in the 1950s and formally
2231 established in his seminal text (Isaacs, 1965). A classic illustrative problem is the **homicidal**
2232 **chauffeur problem** (Isaacs, 1965; Merz, 1974; Patsko and Turova, 2001), which models a pursuit-
2233 evasion scenario between a fast but unmaneuverable pursuer and a slow but agile evader. Today,
2234 differential games are widely applied to solve complex continuous-time problems in economics,
2235 military strategy, and engineering (Dockner et al., 2000).

2236 6.1.2.1 Continuous-time POSG formulation

2237 The POSG formulation of Section 3.5 generalizes to continuous time by replacing the discrete
2238 transition kernel T with a general **jump-diffusion stochastic ODE**:

$$2239 \quad ds_t = \mu(s_t, a_t) dt + \sigma(s_t, a_t) dW_t + \int_{z \sim N(dt)} h(s_t, a_t, z) \quad (6.1)$$

2239 whose components are as follows:

- 2240 • t is the time index.
- 2241 • s_t is the instantaneous state.
- 2242 • a_t is the instantaneous joint action.
- 2243 • W_t is a standard Wiener process (Wiener, 1923) (i.e., Brownian motion), representing continuous,
2244 normally distributed noise.
- 2245 • z is a jump size for the state.
- 2246 • $N(dt)$ is a Poisson random measure (Poisson, 1837; Wiener, 1938), which is the distribution of
2247 jump sizes that occur in the infinitesimal time interval dt .
- 2248 • μ is the drift coefficient, representing the deterministic rate of change.
- 2249 • σ is the diffusion coefficient (or volatility), scaling the continuous random noise.
- 2250 • h is the jump amplitude function.

2251 The return is the Itô integral (Itô, 1944) of the instantaneous rewards r_t weighted by the (possibly
2252 matrix-valued) discounts γ_t .

2253 6.1.2.2 Proposed approach

2254 The existing components of our framework generalize as follows.

- 2255 1. **Trajectory rollout:** Episodes are sampled by numerically integrating the SDE over the time
2256 horizon. In smooth regions of the state space, we use the Dormand–Prince adaptive-stepsize
2257 method (Dormand and Prince, 1980; Hairer, Nørsett, and Wanner, 1993), which automatically
2258 shrinks the stepsize to maintain a target truncation error. Near optimal switching surfaces—
2259 points at which a player’s best action changes discontinuously—we use event detection to locate
2260 the switch time precisely and restart the integrator on either side, preventing the solver from
2261 smearing the discontinuity. We will use JAX-native ODE/SDE solvers like Diffrax (Kidger,
2262 2021), so that the entire rollout remains a pure function of the P2SN parameters and a PRNG
2263 (pseudo-random number generator) key, consistent with our existing JPSPG implementation.
- 2264 2. **Parameter updates:** We apply the JPSPG estimator (Section 5.3.1) to the composition
2265 (P2SN parameters) \rightarrow (integrated trajectory) \rightarrow (episode return). This sidesteps the two
2266 standard obstacles to differentiating through an SDE solver: discontinuities at switching
2267 surfaces, where the adjoint method fails without careful treatment, and the memory cost of
2268 reverse-mode differentiation over long horizons. The constant per-iteration evaluation count of
2269 JPSPG in the number of players is essential here, since continuous-time rollouts are individually
2270 expensive.
- 2271 3. **Best-response estimation:** Exploitability is estimated, as in Section 6.1.1, by training a
2272 single-agent best response against the frozen population policy, which is now a continuous-time
2273 stochastic optimal control problem. If the environment is smooth over time, we can train a
2274 continuous-time actor against the frozen opponents’ P2SN queries. If it is not (such as in
2275 discrete collision events), we fall back to single-player JPSPG on the rollout return.

2276 6.1.2.3 Open questions

2277 These require experimental resolution:

- 2278 1. **Variance–tolerance coupling:** JPSPG’s variance scales with the perturbation radius σ .
2279 Dormand–Prince’s local truncation error scales with the step tolerance ε . The interaction of
2280 the two has not been studied, to our knowledge. We will characterize it empirically and, if
2281 needed, couple the two hyperparameters.
- 2282 2. **Convergence criterion:** Unlike the discrete-time case, exploitability is itself estimated via
2283 an inner optimization that never terminates cleanly. We will report exploitability trajectories
2284 alongside an explicit compute budget, following the convention of our earlier work (Section 5).
- 2285 3. **Failure mode:** If the approach is infeasible on a benchmark (e.g., if best-response training
2286 does not converge within budget), we fall back to a fine-grained discrete approximation and
2287 report the gap.

2288 6.1.2.4 Prior work

2289 In a differential game, players continuously make decisions to optimize their individual payoff functions
2290 over a specified time horizon. A well-known category is the **two-player zero-sum differential**

2291 **game**, where one player’s gain is exactly the other player’s loss (Başar and Olsder, 1999). The
 2292 fundamental mathematical tool for solving such zero-sum games is the **Hamilton–Jacobi–Isaacs**
 2293 (HJI) equation. The continuous-time HJI equation is

$$\min_u \max_v [\nabla V \cdot f(x, u, v) + L(x, u, v)] = 0 \quad (6.2)$$

2294 where V is the value function, u and v are the control inputs, and f represents the system dynamics
 2295 (Friedman, 1971b). The HJI equation characterizes equilibrium only in the two-player zero-sum case.
 2296 For general-sum and many-player settings, our success metric remains empirical exploitability, not
 2297 HJI residual.

2298 Solving differential games analytically is rarely feasible, requiring **numerical integration**
 2299 **techniques** to compute the optimal trajectories and value functions (Bardi and Capuzzo-Dolcetta,
 2300 1997). One approach transforms the PDE into a system of ODEs, using the method of characteristics.
 2301 These characteristic ODEs form a two-point boundary value problem that can be solved using
 2302 iterative numerical integration methods like shooting methods or multiple shooting. Standard
 2303 explicit methods, such as the classic fixed-step fourth-order Runge–Kutta method, often struggle
 2304 with the discontinuities and switching surfaces inherent in optimal control problems. To efficiently
 2305 resolve these sharp gradients and optimal switching times, practitioners employ **adaptive stepsize**
 2306 **algorithms**, such as the Runge–Kutta–Fehlberg or Dormand–Prince methods (Hairer, Nørsett, and
 2307 Wanner, 1993). These adaptive integrators compute two simultaneous approximations of different
 2308 orders to estimate the local truncation error at each time step. If the estimated error exceeds a
 2309 predefined tolerance, the solver automatically reduces the stepsize to maintain accuracy, whereas it
 2310 increases the stepsize in smooth regions to conserve computational resources.

2311 6.1.3 Using ApproxED-BRF to train the P2SN

2312 Our original P2SN paper uses SPSG (Section 5.4.1), the analogue of simultaneous gradient ascent
 2313 (Section 4.3) for the shared-parameter case, to evolve the parameters in time. However, as we
 2314 suggested in the ApproxED paper (Section 5.2.1.1), we could use an alternative parameter dynamics
 2315 based on **minimization of a local approximation of exploitability**. Specifically, we could use
 2316 ApproxED-BRF. This could lead to convergence to equilibrium in a wider range of cases than the
 2317 simultaneous gradient ascent dynamics, as we observed in the ApproxED paper.

2318 Now, in our case, the strategy profile is represented by a P2SN. Therefore, the BRF has to map
 2319 this P2SN to another P2SN. Since the strategy profile is itself a neural network, this means the BRF
 2320 must be a **hypernetwork** (Schmidhuber, 1992; Ha, Dai, and Le, 2017; Lorraine and Duvenaud,
 2321 2018; MacKay et al., 2019; Bae and Grosse, 2020), that is, a neural network that maps one neural
 2322 network to another. In our ApproxED paper, we already used hypernetworks for strategy profiles in
 2323 which the strategies were policy networks.

2324 Recall from Section 5.4.1 that the shared-parameter simultaneous gradient (SPSG) is defined as

$$\left[\frac{d}{d\phi} \int_{i \sim \mu} u(s_\theta[i \mapsto (s_\phi)_i])_i \right]_{\phi=\theta} \quad (6.3)$$

2325 Recall from Section 5.2.1.1 that the BRF evolves as follows:

$$\dot{x} = -\nabla_x \text{NI}(x, b_\theta(x)) \quad (6.4)$$

$$\dot{\theta} = +\nabla_\theta \text{NI}(x, b_\theta(x)) \quad (6.5)$$

2326 If we let θ refer to the parameters of the P2SN and ϕ refer to the parameters of the best-response
 2327 hypernetwork (which takes the P2SN’s parameters as input), we have

$$\dot{\theta} = -\nabla_{\theta} \int_{i \sim \mu} (u(s_{\theta}[i \mapsto b_{\phi}(\theta)_i])_i - u(s_{\theta})_i) \quad (6.6)$$

$$\dot{\phi} = +\nabla_{\phi} \int_{i \sim \mu} (u(s_{\theta}[i \mapsto b_{\phi}(\theta)_i])_i - u(s_{\theta})_i) \quad (6.7)$$

2328 6.1.4 Mitigating catastrophic forgetting in ApproxED-BRF

2329 Recall that ApproxED-BRF (Section 5.2.1.1) trains a best response function against the current
 2330 strategy profile. The motivation is that a function can model best responses for many strategy
 2331 profiles simultaneously, and thus in principle can “retain” memory of what good responses are
 2332 for prior strategy profiles, avoiding the cycling problems encountered by naive learning, that is,
 2333 simultaneous gradient ascent (Section 4.3). However, since the BRF trains only against the current
 2334 strategy profile on each iteration, it is in principle still possible for the BRF to “forget” what it has
 2335 learned in the distant past, i.e., how to best respond to strategy profiles that were only encountered
 2336 in the distant past.

2337 In the ML literature, **catastrophic forgetting** (or catastrophic interference) is a fundamental
 2338 limitation in artificial neural networks where the model completely and abruptly forgets previously
 2339 learned information upon learning new data or tasks (McCloskey and Cohen, 1989). This occurs
 2340 because the weights (parameters) of the network that were optimized for a previous task are
 2341 heavily altered through backpropagation to minimize the error on the new task. As a result, the
 2342 network’s performance on the old task drastically degrades (French, 1999). This prevents standard
 2343 neural networks from achieving “continual” or “lifelong” learning, a key feature of human biological
 2344 intelligence.

2345 Researchers have proposed various strategies to mitigate catastrophic forgetting, which generally
 2346 fall into three main categories.

2347 **Regularization-based methods.** These methods introduce an additional term to the loss function
 2348 to penalize changes to weights that are deemed critical for previously learned tasks.

- 2349 • **Elastic weight consolidation (EWC):** EWC slows down learning on certain weights based
 2350 on how important they are to previously seen tasks, estimating this importance using the Fisher
 2351 information matrix (Kirkpatrick et al., 2017).
- 2352 • **Learning without forgetting (LwF):** This approach utilizes knowledge distillation. It uses
 2353 the model’s own prior outputs on new task data to create a regularization target, ensuring the
 2354 network’s responses to old tasks do not drift significantly (Li and Hoiem, 2017).

2355 **Architectural and parameter isolation methods.** These methods dynamically alter the
 2356 architecture of the neural network to dedicate distinct subsets of parameters to different tasks.

- 2357 • **Progressive neural networks:** This method instantiates a new neural network column for each
 2358 new task. It prevents forgetting by freezing the weights of previous columns and only allowing
 2359 lateral connections to the new column (Rusu et al., 2016).

2360 **Rehearsal and replay methods.** Rehearsal techniques prevent forgetting by periodically re-
 2361 training the network on a subset of the old data while learning the new data.

- 2362 • **Experience replay:** A memory buffer stores exact examples from previous tasks, which are
 2363 interleaved with new data during training (Rebuffi et al., 2017).
- 2364 • **Generative replay:** Instead of storing raw data, a separate generative model (like a GAN)
 2365 is trained to generate pseudo-data from past distributions, which is then replayed to the main
 2366 network (Shin et al., 2017).

2367 We opt for the replay approach. That is, instead of relying only on the implicit flexibility of the
 2368 BRF, we go one step further and *explicitly* train the BRF against multiple strategy profiles on each
 2369 iteration. Specifically, instead of training against only the current strategy profile (i.e., the P2SN’s
 2370 parameters), we maintain an **experience replay buffer** of past strategy profiles.

2371 We can fill the replay buffer via **reservoir sampling** (Vitter, 1985). Reservoir sampling is a
 2372 family of randomized algorithms used to select a simple random sample of k items from a population
 2373 of unknown or vast size n , where n is either too large to fit into memory or is a continuous data
 2374 stream. It allows one to maintain a representative sample at any point during the stream.

2375 By utilizing reservoir sampling to manage the experience replay buffer, continual learning models
 2376 achieve the following benefits.

- 2377 • **Uniform historical representation:** Reservoir sampling mathematically guarantees that at
 2378 any given time step n , every single data point seen so far has an exact equal probability (k/n) of
 2379 being in the buffer. This ensures the replay buffer remains an unbiased, uniform sample of the
 2380 entire data distribution seen over the model’s lifetime.
- 2381 • **Memory efficiency:** The model does not need to store the entire dataset or even know the
 2382 total size of the dataset in advance. It only requires the k slots in the buffer and a counter for n ,
 2383 making it highly efficient for continuous learning environments.

2384 In summary, we store past strategy profiles at random, and on every iteration, we sample one
 2385 and train the BRF against it. This makes the BRF less likely to overfit to the current strategy
 2386 profile (or recent profiles), and mitigates catastrophic forgetting.

2387 6.1.5 Modifying BRF’s update scheme to use other learning dynamics

2388 BRF (Section 5.2.1.1) resembles simultaneous gradient ascent (Section 4.3) applied at the level
 2389 of strategy profiles rather than individual strategies. Specifically, recall from Section 4.3 that
 2390 simultaneous gradient ascent is as follows for a two-player zero-sum game:

$$\dot{x} = -\nabla_x u(x, y) \tag{6.8}$$

$$\dot{y} = +\nabla_y u(x, y) \tag{6.9}$$

2391 Recall from Section 5.2.1.1 that the BRF evolves as follows:

$$\dot{x} = -\nabla_x \text{NI}(x, b_\theta(x)) \tag{6.10}$$

$$\dot{\theta} = +\nabla_\theta \text{NI}(x, b_\theta(x)) \tag{6.11}$$

2392 This is because the goal is to find

$$x^* \in \underset{x \in \mathcal{X}}{\text{argmin}} \text{NI}(x, b(x)) \tag{6.12}$$

2393 such that $b : \mathcal{X} \rightarrow \mathcal{Y}$ is a function that satisfies

$$b(x) \in \operatorname{argmax}_{y \in \mathcal{Y}} \text{NI}(x, y) \quad (6.13)$$

2394 In Section 4.3, we saw that there are other learning dynamics that modify the simple simultaneous
 2395 gradient ascent scheme. These include extragradient and optimistic gradient, as well as other methods
 2396 that involve higher-order derivatives. We could substitute these for the simultaneous gradient ascent
 2397 that BRF uses at the level of strategy profiles (versus the BRF’s parameters).

2398 Since the game that ApproxED-BRF is trying to solve is technically a Stackelberg game (we
 2399 are trying to find the min-max specifically, not the max-min, so the order matters), we can also
 2400 consider dynamics made specifically for **Stackelberg games**. Examples of these can be found in
 2401 Fiez, Chasnov, and Ratliff (2019), Fiez, Chasnov, and Ratliff (2020), Goktas and Greenwald (2021),
 2402 Goktas and Greenwald (2022b), Goktas, Zhao, and Greenwald (2022b), and Goktas, Zhao, and
 2403 Greenwald (2022a).

2404 6.1.6 Joint-perturbation estimator for other learning dynamics

2405 For this objective, the goal is to develop analogues of the JPSPG estimator (Section 5.3.1) for
 2406 other learning dynamics (Section 4.3), including ones that involve higher-order derivatives and
 2407 Hessian-vector products, so that these can be used efficiently as well.

2408 Recall from Section 4.3 that the simultaneous gradient is

$$g_i = \frac{du(x)_i}{dx_i} \quad (6.14)$$

2409 Recall from Section 5.3.1 that the simultaneous pseudo-gradient estimator is

$$\hat{g}_i = \frac{1}{\sigma} u(x[i \mapsto x_i + \sigma z_i])_i z_i \quad (6.15)$$

2410 Recall from Section 5.3.1 that the joint-perturbation JPSPG estimator is

$$\hat{g} = \frac{1}{\sigma} u(x + \sigma z) \odot z \quad (6.16)$$

2411 This estimator can be re-used for both extragradient and optimistic gradient, which are comprised
 2412 simply of expressions involving the simultaneous gradient. For the dynamics involving higher-order
 2413 derivatives, the answer is less trivial. For these, the formulation of each dynamics in terms of
 2414 the stop-gradient operator, as documented in Letcher (2018), is helpful. The expression inside
 2415 each stop-gradient does not undergo perturbation for the derivative that is being computed via
 2416 perturbation. For illustration, we consider a few of the dynamics. The same process can be carried
 2417 out for the other dynamics in an analogous way.

2418 To construct the updates elegantly, we define two base estimators using a single joint perturbation
 2419 evaluation. The first-order pseudo-gradient for player i uses the antithetic central difference:

$$\hat{g}_i = \frac{u(x + \sigma z)_i - u(x - \sigma z)_i}{2\sigma} z_i \quad (6.17)$$

2420 The scalar factor \hat{d}_i captures the **second-order curvature** of the utility function along the
 2421 specific perturbation vector, via a second-order central difference stencil:

$$\hat{d}_i = \frac{u(x + \sigma z)_i + u(x - \sigma z)_i - 2u(x)_i}{\sigma^2} \quad (6.18)$$

2422 By leveraging the outer product of the perturbations, the cross-agent Jacobian blocks simplify
 2423 to $\hat{J}_{ij} = \hat{d}_i z_i z_j^T$. This allows us to collapse the complex matrix operations of LOLA, CO, and SGA
 2424 into isolated scalar products multiplied by the base perturbation vector z_i .

2425 For the continuous-time differential equations, we define the exact individual gradient as $g_i \triangleq$
 2426 $\nabla_{x_i} u(x)_i$.

2427 Simultaneous gradient ascent simply evolves according to the simultaneous gradient:

$$\dot{x}_i = g_i \tag{6.19}$$

$$x_i \leftarrow x_i + \alpha_i \hat{g}_i \tag{6.20}$$

2428 LOLA shapes the anticipated naive updates of opponents based on their individual learning
 2429 rates η_j :

$$\dot{x}_i = g_i + \sum_{j \neq i} \eta_j (\nabla_{x_j} g_i) g_j \tag{6.21}$$

$$x_i \leftarrow x_i + \alpha_i \left(\hat{g}_i + \hat{d}_i \left(\sum_{j \neq i} \eta_j z_j^T \hat{g}_j \right) z_i \right) \tag{6.22}$$

2430 CO stabilizes the joint vector field by penalizing it with the transposed Jacobian-vector product
 2431 scaled by a hyperparameter γ :

$$\dot{x}_i = g_i - \gamma \sum_j (\nabla_{x_i} g_j) g_j \tag{6.23}$$

$$x_i \leftarrow x_i + \alpha_i \left(\hat{g}_i - \gamma \left(\sum_j \hat{d}_j z_j^T \hat{g}_j \right) z_i \right) \tag{6.24}$$

2432 SGA isolates the rotational dynamics by applying the antisymmetric part of the game Jacobian
 2433 scaled by λ :

$$\dot{x}_i = g_i + \frac{\lambda}{2} \sum_j (\nabla_{x_j} g_i - \nabla_{x_i} g_j) g_j \tag{6.25}$$

$$x_i \leftarrow x_i + \alpha_i \left(\hat{g}_i + \frac{\lambda}{2} \left(\hat{d}_i \sum_j z_j^T \hat{g}_j - \sum_j \hat{d}_j z_j^T \hat{g}_j \right) z_i \right) \tag{6.26}$$

2434 By extracting the scalar finite-difference terms upfront, these update rules collapse the cross-agent
 2435 interactions into simple scalar multipliers applied to the perturbation vector z_i . This results in
 2436 methods with memory and computational complexity that is strictly linear with respect to the
 2437 parameter count, while maintaining the constant function evaluation complexity of JPSPG.

2438 6.2 Benchmarks

2439 In this section, we describe our proposed benchmarks. These environments naturally exhibit the
 2440 complexities outlined in Section 1.2, including infinite heterogeneous populations, matrix-valued

2441 discounting, and discontinuous payoffs. We describe each environment as a POSG, as defined in
 2442 Section 3.5. Across all benchmarks, we assume I is a non-atomic measure space (for example, the
 2443 unit interval $[0, 1]$ or unit square $[0, 1]^2$ equipped with the Lebesgue measure), modeling an infinite
 2444 population of players. Our goal is to show that, for each of these environments, our method converges
 2445 to Nash equilibrium, as measured by approximate exploitability.

2446 6.2.1 Epidemic: Epidemiological contagion and pandemic mitigation

2447 Drawing from the behavioral epidemiology literature (Reluga, 2010; Acemoglu et al., 2020), this is a
 2448 mean-field game of viral contagion. Each individual $i \in I$ has a private state $x_i = (v_i, r_i) \in \mathbb{R}^2 = X$
 2449 consisting of a viral load v_i and financial reserve r_i . The global state is X^I , consisting of the
 2450 individuals' private states. The components of the POSG are as follows:

- 2451 • $I = [0, 1]$ is the non-atomic continuum of individuals.
- 2452 • $S = X^I$ is the global state space.
- 2453 • $O_i = X \times \mathbb{R}$ contains the private state x_i alongside a scalar public health signal.
- 2454 • $A_i = [0, 1]$ is the individual's daily socialization effort.
- 2455 • $Z : S \rightarrow \Delta\Pi O$ is the observation kernel. It provides agent i with $o_i = (x_i, \hat{v})$, where $\hat{v} \sim \mathcal{N}(\bar{v}, \sigma^2)$
 2456 is a noisy broadcast of the mean viral load $\bar{v} = \int_{i \sim \mu} v_i$.
- 2457 • $T : S \times \Pi A \rightarrow \Delta S$ is the transition kernel. The viral load v_i evolves stochastically based on a_i
 2458 and the societal infection density derived from ν . Financial reserves update deterministically as
 2459 $r'_i = f_i + w_i a_i - c_i$, where w_i and c_i denote cohort-specific wages and living costs, respectively.
- 2460 • $R : S \times \Pi A \rightarrow V$ has a strict dual-discontinuity penalty. The utility R_i suffers a large penalty if
 2461 the macroscopic mass of severe cases $\int_{i \sim \mu} \mathbb{1}[v_i > v_{\text{crit}}]$ strictly exceeds a healthcare capacity C , or
 2462 if the agent faces bankruptcy ($r_i \leq 0$).
- 2463 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ is an integral operator defining demographic altruism, with a kernel
 2464 $K(i, j)$ that specifies how $i \in I$ discounts the future returns of $j \in I$.
- 2465 • $\rho \in \Delta S$ initializes the global state to reflect a localized outbreak.

2466 6.2.2 Finance: Financial contagion and systemic risk

2467 Inspired by models of systemic risk (Eisenberg and Noe, 2001; Elliott, Golub, and Jackson, 2014),
 2468 this game models a global financial network. Let $X = \mathbb{R}^k$ represent the private balance sheet of
 2469 institution i , capturing $k \in \mathbb{N}$ distinct balance sheet variables, including capital reserves c_i and
 2470 toxic asset holdings w_i . To capture exogenous market forces, let $M = \mathbb{R}^m$ represent $m \in \mathbb{N}$ global
 2471 macroeconomic indices. The components of the POSG are as follows:

- 2472 • $I = [0, 1]$ is the continuum of financial institutions.
- 2473 • $S = X^I \times M$ couples the microscopic balance sheets with the independent macroeconomic
 2474 indicators.
- 2475 • $O_i = X \times M$ reflects the fact that an institution observes its own balance sheet (one copy of X)
 2476 and the global indices M , but not the balance sheets of other institutions.
- 2477 • $A_i \in \mathbb{R}^d$ denotes the continuous capital allocation vectors for interbank lending and liquidation.

- 2478 • $Z : S \rightarrow \Delta\Pi O$ exposes private state components of i to i .
- 2479 • $T : S \times \Pi A \rightarrow \Delta S$ is a **jump-diffusion process** updating both balance sheets and global indices.
- 2480 Crucially, capital defaults trigger discontinuous cascading shocks across neighboring transition
- 2481 probabilities.
- 2482 • $R : S \times \Pi A \rightarrow V$ imposes regulatory discontinuities. If capital reserves c_i fall below a statutory
- 2483 threshold θ , R_i evaluates strictly to zero and a forced liquidation is triggered.
- 2484 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ acts as a linear operator that dynamically redistributes discounted
- 2485 returns in proportion to institutional equity ownership percentages.
- 2486 • $\rho \in \Delta S$ initializes the market with a baseline distribution of interconnected assets and capital.

2487 6.2.3 Grid: Wholesale forward capacity markets

2488 Drawing from mean-field models of smart grids (Borenstein, Bushnell, and Knittel, 1999; Kizilkale,
 2489 Salhab, and Malhamé, 2019), this game models a dynamic electricity market. Let $X = \mathbb{R}_{\geq 0}$ represent
 2490 the physically constrained megawatt capacity c_i of generator i . The components of the POSG are as
 2491 follows:

- 2492 • $I = [0, 1]$ is the population of heterogeneous power generators.
- 2493 • $S = X^I \times \mathbb{R}_{\geq 0}$ couples the microscopic capacities with a global grid frequency.
- 2494 • $O_i = X \times \mathbb{R}_{\geq 0}$ provides local capacity measurements alongside a stochastically noisy reading of
- 2495 grid frequency.
- 2496 • $A_i \in \mathbb{R}_{\geq 0}^2$ represents the continuous parameters of a step-function supply bid, comprising an
- 2497 offered volume $v_i \leq c_i$ and a minimum acceptable price p_i . These yield the individual supply
- 2498 curve $S_i(P) = v_i \mathbb{1}[P \geq p_i]$.
- 2499 • $Z : S \rightarrow \Delta\Pi O$ obfuscates the global supply curve and competing bids.
- 2500 • $T : S \times \Pi A \rightarrow \Delta S$ evolves individual capacities c_i via mean-reverting **Ornstein-Uhlenbeck**
- 2501 **processes** representing weather stochasticity, while the grid frequency updates reactively based
- 2502 on aggregate market bid volumes.
- 2503 • $R : S \times \mathbf{A} \rightarrow V$ imposes a strict market-clearing discontinuity via a uniform-price auction. Let
- 2504 $D \in \mathbb{R}_{\geq 0}$ represent the macroscopic inelastic electricity demand. The environment computes
- 2505 the aggregate supply curve by integrating the step-function bids across the continuum: $S(P) =$
- 2506 $\int_{j \sim I} v_j \mathbb{1}(P \geq p_j)$. The market establishes the uniform clearing price P^* at the exact intersection
- 2507 of supply and demand:

$$P^* = \inf\{P \in \mathbb{R}_{\geq 0} : S(P) \geq D\} \quad (6.27)$$

2508 The reward R_i for generator i is then evaluated as a piecewise function, creating a sharp
 2509 discontinuity for bids exceeding the clearing price:

$$R_i(s, \mathbf{a}) = \begin{cases} v_i P^* & \text{if } p_i \leq P^* \\ 0 & \text{if } p_i > P^* \end{cases} \quad (6.28)$$

- 2510 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ encodes regional stability pacts, internalizing the grid stability of
- 2511 physically interconnected neighbors.
- 2512 • $\rho \in \Delta S$ samples the baseline weather states and initial generator capacities.

2513 6.2.4 Foraging: Evolutionary foraging with kin selection

2514 Drawing from the foundations of evolutionary game theory (Hamilton, 1964; McNamara and Houston,
 2515 1996), we formulate a spatial foraging game. Let $X = \mathbb{R}_{\geq 0} \times \mathbb{R}^2$ define the private biological state,
 2516 comprising internal energy reserves e_i and 2D spatial coordinates p_i . Let $E = L^2(\mathbb{R}^2, \mathbb{R}_{\geq 0})$ be the
 2517 independent, continuous spatial density field of regenerating nutrients. The components of the POSG
 2518 are as follows:

- 2519 • $I = [0, 1]$ is the continuum of competing biological organisms.
- 2520 • $S = X^I \times E$ combines the microscopic states of organisms with the global nutrient field.
- 2521 • $O_i = \mathbb{R}_{\geq 0} \times \mathbb{R}_{\geq 0}$ equips the agent with knowledge of its internal energy e_i and a localized sensory
 2522 integral of nearby nutrients centered at p_i .
- 2523 • $A_i \in \mathbb{R}^2$ is the continuous velocity vector dictating spatial movement.
- 2524 • $Z : S \rightarrow \Delta\Pi O$ restricts spatial observation to an area around p_i .
- 2525 • $T : S \times \Pi A \rightarrow \Delta S$ updates spatial coordinates $p'_i = p_i + a_i$, depletes energy $e'_i = e_i - c\|a_i\|$, and
 2526 stochastically regenerates the encompassing nutrient field E .
- 2527 • $R : S \times \Pi A \rightarrow V$ enforces a biological starvation discontinuity. Utility scales with successful
 2528 nutrient intake, but becomes permanently zero if the internal energy e_i falls below a critical
 2529 survival threshold e_{crit} .
- 2530 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ weighs the returns of agent j in agent i 's utility by genetic relatedness r_{ij} ,
 2531 realizing **Hamilton's rule** of kin selection¹, weighing the returns of genetically-similar agents.
- 2532 • $\rho \in \Delta S$ samples the initial spatial distribution of both agents and available nutrients.

2533 6.2.5 Orbit: Orbital constellation management and debris avoidance

2534 Inspired by the physical challenges of low Earth orbit mega-constellations (Kessler and Cour-Palais,
 2535 1978; Alfriend et al., 2009), we model an orbital routing game. Let $X = \mathbb{R}^6 \times \mathbb{R}_{\geq 0}$ represent the
 2536 private state, capturing the 6D orbital vector q_i (position and velocity) alongside finite chemical fuel
 2537 reserves f_i . The components of the POSG are as follows:

- 2538 • $I = [0, 1]$ represents the continuum of active satellites interspersed with untrackable micro-debris.
- 2539 • $S = X^I$ distills the global state entirely to the microscopic profile of all orbital objects.
- 2540 • $O_i = X \times \mathbb{R}_{\geq 0}$ provides the exact internal state and noisy radar estimations of proximate objects.
- 2541 • $A_i \in \mathbb{R}^3$ dictates the continuous thrust vector allocated for station-keeping maneuvers.
- 2542 • $Z : S \rightarrow \Delta\Pi O$ simulates radar noise, sensor degradation, and periodic line-of-sight signal
 2543 occlusion.

¹Hamilton's rule of kin selection (Hamilton, 1964) is a mathematical formula that explains why an animal might act altruistically, sacrificing its own fitness to help another, even when it seems to go against its own survival. It suggests that "altruistic" genes can evolve if the benefit to a relative is high enough to outweigh the cost to the individual. The rule is that genes for a particular behavior should increase in frequency when $rB > C$, where r is the genetic relatedness of the recipient to the actor, B is the additional reproductive benefit gained by the recipient of the altruistic act, and C is the reproductive cost to the individual performing the act.

- 2544 • $T : S \times \Pi A \rightarrow \Delta S$ integrates non-linear orbital mechanics to constantly update q_i , while strictly
2545 depleting fuel as $f'_i = f_i - c\|a_i\|$, where c represents the specific fuel consumption (or the inverse
2546 of the propulsion system’s efficiency).
- 2547 • $R : S \times \Pi A \rightarrow V$ models the Kessler syndrome². The utility R_i evaluates to a massive penalty
2548 if the spatial distance between q_i and any q_j drops below a physical collision radius, or if the
2549 satellite burns through all remaining fuel ($f_i \leq 0$).
- 2550 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ represents overarching geopolitical alliances, assigning high positive
2551 weights to allied orbital assets and non-positive weights to adversaries.
- 2552 • $\rho \in \Delta S$ initializes orbits derived directly from empirical aerospace tracking catalogs.

2553 6.2.6 Traffic: Macroscopic urban traffic and fleet routing

2554 Drawing from mean-field routing games (Wardrop, 1952; Bressan, 2015), this game models the
2555 fluid dynamics of continuous urban traffic. Let $X = \mathbb{R}^2 \times \mathbb{R}^2 \times \mathbb{R}_{\geq 0}$ represent the state of driver i ,
2556 structurally containing position p_i , velocity v_i , and an elapsed trip time t_i . The empirical spatial
2557 density field ν is derived from the profile of all driver positions. The components of the POSG are
2558 as follows:

- 2559 • $I = [0, 1]$ is the continuum of vehicles navigating the road network.
- 2560 • $S = X^I$ maintains the global state strictly through the kinematics and trip times of all individual
2561 agents.
- 2562 • $O_i = X \times \mathbb{R}_{\geq 0}$ is an observation of local kinematics alongside a noisy local measurement of ν .
- 2563 • $A_i \in \mathbb{R}^2$ represents the combined acceleration and steering vector.
- 2564 • $Z : S \rightarrow \Delta \Pi O$ limits observation to the local density.
- 2565 • $T : S \times \Pi A \rightarrow \Delta S$ updates p_i and v_i subject to a non-linear drag generated directly by the local
2566 density derived from ν . Elapsed time t_i increments deterministically.
- 2567 • $R : S \times \Pi A \rightarrow V$ implements gridlock and deadline discontinuities. The utility R_i monotonically
2568 decays with trip time t_i , but drops if local density breaches a maximum flow threshold or if t_i
2569 violates a strict delivery deadline.
- 2570 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ enables ride-sharing fleet coordination, distributing discounted returns
2571 across vehicles operating within the same corporate network, which have an incentive to optimize
2572 *global* systemic throughput.
- 2573 • $\rho \in \Delta S$ initializes the macroscopic spatial profile of vehicles by sampling from historically-
2574 grounded spatial gravity and bottleneck models (Voorhees, 1955; Wilson, 1967; Vickrey, 1969;
2575 Arnott, De Palma, and Lindsey, 1990), generating realistic, asymmetrical origin-destination
2576 demand across the density field.

²Kessler syndrome (Kessler and Cour-Palais, 1978) describes a situation in which the density of objects in low Earth orbit (LEO) becomes so high due to space pollution that collisions between these objects cascade, exponentially increasing the amount of space debris over time.

2577 6.3 Codebase

2578 Our codebase is written in **Python** (Rossum and Drake, 1995) and uses the following libraries:

- 2579 • **JAX** (Bradbury et al., 2018): a library for high-performance numerical computing and large-scale
2580 machine learning.
- 2581 • **Flax** (Heek et al., 2023): a neural network library and ecosystem for JAX designed for flexibility.
- 2582 • **Optax** (DeepMind et al., 2020): a gradient processing and optimization library for JAX. It
2583 contains implementations of the most common optimizers used in machine learning.
- 2584 • **DiffraX** (Kidger, 2021): a JAX-based library providing numerical differential equation solvers
2585 that is autodifferentiable and GPU-capable.
- 2586 • **Mctx** (DeepMind et al., 2020): JAX-native implementation of Monte Carlo tree search (MCTS)
2587 algorithms.
- 2588 • **Matplotlib** (Hunter, 2007): a comprehensive library for creating static, animated, and interactive
2589 visualizations.
- 2590 • **SciPy** (Virtanen et al., 2020): an open-source software for mathematics, science, and engineering.
2591 It includes modules for statistics, optimization, integration, linear algebra, Fourier transforms,
2592 signal and image processing, ODE solvers, and more.
- 2593 • **NumPy** (Harris et al., 2020): the fundamental package for scientific computing with Python.
- 2594 • **Pandas** (The Pandas development team, 2020; McKinney, 2010): a flexible and powerful data
2595 analysis / manipulation library for Python, providing labeled data structures.
- 2596 • **PyCairo**: Python bindings for the Cairo graphics library (Worth and Packard, 2003).
- 2597 • **CVXPY** (Diamond and Boyd, 2016; Agrawal et al., 2018): a Python-embedded modeling
2598 language for convex optimization problems.
- 2599 • **OTT-JAX** (Cuturi et al., 2022): JAX package providing tools for solving optimal transport
2600 problems (e.g., Wasserstein distance between point clouds).
- 2601 • **h5py**: a thin, Pythonic wrapper around HDF5 (Koranne, 2010; The HDF Group, 2024), a
2602 high-performance data management and storage suite.

2603 I have extensive experience with Python, JAX, and all of these libraries. I use them daily to code
2604 my environments, agents, and learning algorithms. For evidence of this, below are my open-source
2605 contributions to JAX and Optax:


- 2606 • <https://github.com/jax-ml/jax/pulls?q=is:pr+author:carlosmartin>
- 2607 • <https://github.com/google-deeppmind/optax/pulls?q=is:pr+author:carlosmartin>

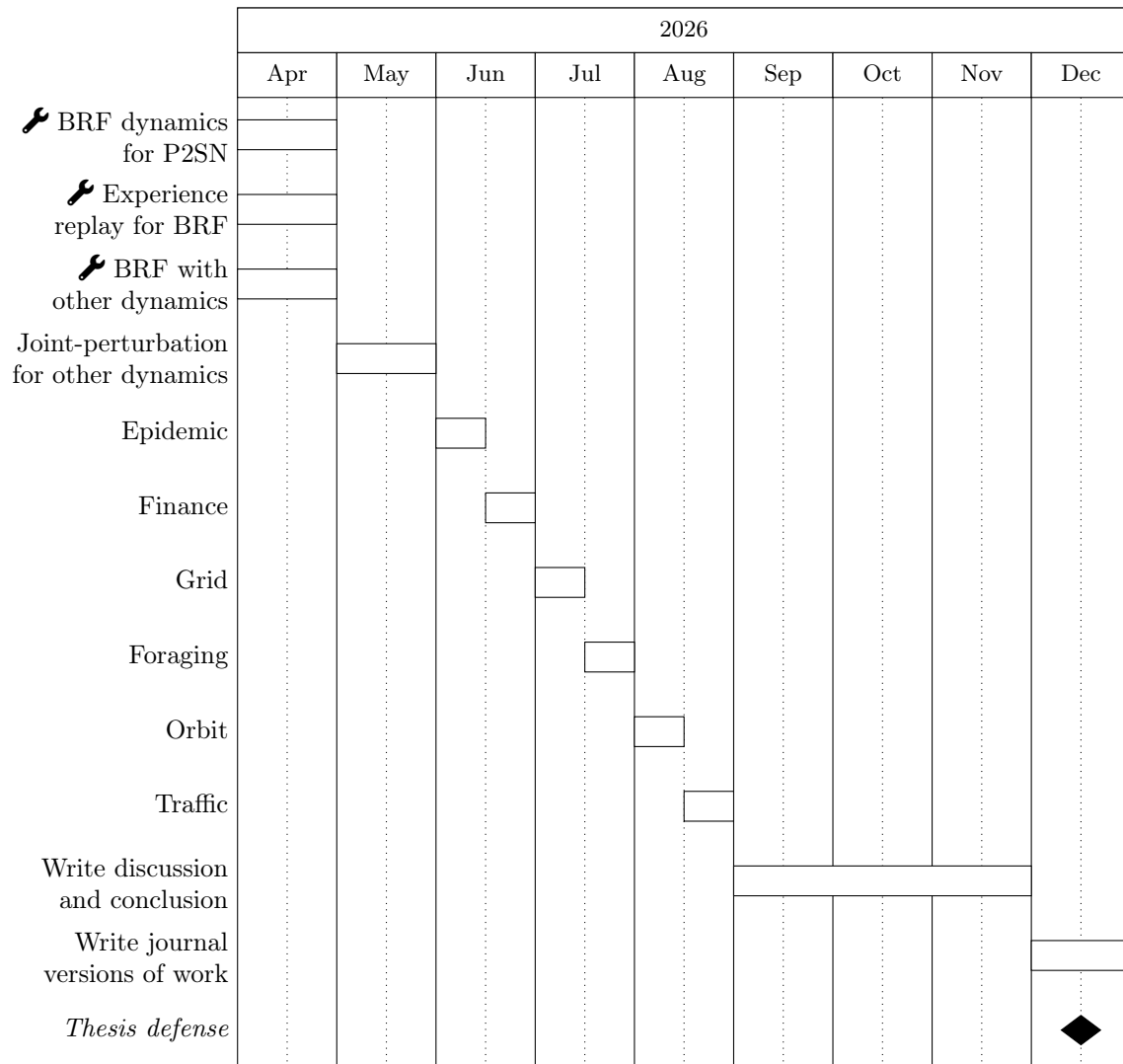
2608 6.4 Available compute

2609 My lab has a dedicated computing cluster, the Marvel computing cluster. It has 1 node with an
2610 AMD EPYC 7742 64-Core Processor and 8 NVIDIA A100 GPUs. It also has 14 other nodes, each
2611 with an AMD EPYC 7282 16-Core Processor. It also has 5 other nodes, each with an AMD Opteron
2612 Processor 6272.

2613 I also have access to the Google Cloud ORCHARD Cluster managed by CMU. Each job can run
 2614 for 12 hours and has access to a maximum of 1 GPU.

2615 6.5 Timeline

2616 My proposed timeline is as follows. As stated before, most of the code and writeup is **already done**.
 2617 Also, experiments are running continuously in the background on my lab's cluster, with occasional
 2618 monitoring, so experiments can overlap with either the programming or the writeup. The symbol
 2619  is used to mark work that is almost complete.



6.6 Stretch goals

The contributions outlined above form the foundational core of this thesis and are sufficient to validate it. However, if time and computational resources permit, I also plan to explore the following optional extensions. While not strictly necessary for the completion of the dissertation, these speculative directions build naturally upon my core framework, and represent potentially interesting avenues for future research.

- Extend to **alternative solution concepts**. These include the following:
 - Correlated equilibria (CE) (Aumann, 1974).
 - Coarse correlated equilibria (CCE) (Moulin and Vial, 1978).
 - Quantal response equilibria (QRE) (McKelvey and Palfrey, 1995; McKelvey and Palfrey, 1998).
 - α -rank (Omidshafiei et al., 2019).
- Extend to **mechanism design**, that is, tuning game parameters so that the resulting equilibria have desirable properties from the perspective of a social planner. There is a vast literature on mechanism design. For this, we could use **bilevel optimization** techniques. The state of the art in continuous bilevel optimization uses implicit differentiation and variance-reduced stochastic approximation to efficiently estimate inverse Hessians and compute hypergradients at scale (Lorraine, Vicol, and Duvenaud, 2020; Ji, Yang, and Liang, 2021). For discrete or mixed-integer settings, state-of-the-art methods use specialized branch-and-bound frameworks applied to Karush–Kuhn–Tucker (KKT) or value-function reformulations (Kleinert et al., 2021).
- Extend to **inverse multiagent reinforcement learning** (IRL). This means inferring characteristics of a game from players’ behavior. I have already done some work along these lines in Martin and Sandholm (2021a). To scale this up to more general games where exact Bayesian inference is intractable, we could use some tractable approximations. For example, we could use the following techniques.
 - Markov chain Monte Carlo (MCMC) (Metropolis et al., 1953; Hastings, 1970) is a class of algorithms that enables sampling from complex, high-dimensional probability distributions by constructing a Markov chain that has the desired distribution as its equilibrium state. These include the Metropolis–Hastings (MH) algorithm (Hastings, 1970; Metropolis et al., 1953), Hamiltonian Monte Carlo (HMC) (Duane et al., 1987; Neal, 1996), and No-U-Turn Sampling (NUTS) (Hoffman and Gelman, 2014).
 - Dropout (Srivastava et al., 2014) is an effective way to represent model uncertainty in deep learning (Gal and Ghahramani, 2016).
 - Likewise, DropConnect (Wan et al., 2013) is effective in modeling uncertainty of Bayesian deep networks (Mobiny et al., 2021).
- Extend to **infinite utilities**. To be more precise, these would be non-Archimedean utilities, under lexicographic ordering or leximin ordering. These could be handled via **multilevel optimization**.
- Incorporate **efficient exploration** techniques, as we did in Martin and Sandholm (2021b). There is a vast literature on efficient exploration for RL and MARL. We could try to generalize it and apply it to our infinite setting.

- 2660 • Incorporate **execution-time planning** into the agents, as we did in Martin, Boutilier, et al.
2661 (2024) and Martin and Sandholm (2025a).

2662 In addition, we include the following environment as a stretch goal.

2663 6.6.1 Auction: Real-time programmatic ad auctions

2664 We adapt continuous real-time bidding (Edelman, Ostrovsky, and Schwarz, 2007; Cai, Ren, et al.,
2665 2017) into a multi-step POSG. Let $X = \mathbb{R}_{\geq 0}^2$ represent the private state of advertiser i , consisting
2666 of a depleting remaining budget b_i and a hidden conversion valuation v_i . The components of the
2667 POSG are as follows:

- 2668 • $I = [0, 1]$ is the continuum of autonomous advertising agents.
- 2669 • $S = X^I$ is the global state consisting of the agents' profile of budgets and valuations.
- 2670 • $O_i = \mathbb{R}_{\geq 0}$ restricts observation, providing the agent only with its current budget b_i while hiding
2671 its true conversion rate and the states of competitors.
- 2672 • $A_i \in \mathbb{R}_{\geq 0}$ is the continuous financial bid submitted for an impression, strictly bounded by the
2673 remaining budget ($a_i \leq b_i$).
- 2674 • $Z : S \rightarrow \Delta\Pi O$ acts as an information filter, obscuring the latent macroscopic demographic
2675 variables D_t that govern the stochastic evolution of the true conversion valuation v_i , forcing the
2676 agent to infer market regimes strictly from delayed auction outcomes.
- 2677 • $T : S \times \Pi A \rightarrow \Delta S$ deterministically subtracts the clearing price from the winning agent's budget
2678 b_i and stochastically evolves latent valuations v_i via a hidden Markov model.
- 2679 • $R : S \times \Pi A \rightarrow V$ models a generalized second-price auction mechanism. The utility yields
2680 $R_i = v_i - \max_{j \neq i} a_j$ if $a_i > \max_{j \neq i} a_j$, and zero otherwise.
- 2681 • $\Gamma : S \times \Pi A \times S \rightarrow \mathcal{L}(V)$ models advertising agency portfolios, coupling the returns of client
2682 brands operating under the same centralized bidding algorithm.
- 2683 • $\rho \in \Delta S$ assigns the initial campaign budgets and baseline valuation distributions.

2684 In programmatic advertising, the true value of an ad impression (the conversion rate) is not static.
2685 It is driven by macroscopic, shifting demographic trends, such as a viral social media trend making
2686 teenagers suddenly buy a specific product, a seasonal purchasing behavior, or a competitor launching
2687 a rival ad campaign. Advertisers do not get to see these underlying demographic shifts directly. They
2688 only see the downstream effects: whether their bid won, and whether that win eventually resulted in
2689 a conversion. The market is constantly shifting beneath their feet, and they must infer those shifts
2690 purely from incomplete, noisy feedback. To formalize this, we can model the demographic shifts as
2691 a latent stochastic process. Let $D_t \in \mathbb{R}^m$ represent a hidden macroscopic vector tracking $m \in \mathbb{N}$
2692 demographic/market regimes at time t . The true, hidden valuation $v_{i,t}$ is no longer a static number,
2693 but a dynamic function of this latent demographic state: $v_{i,t} = f_i(D_t) + \varepsilon$, where ε is a noise term.

Bibliography

- 2695 Acemoglu, Daron, Victor Chernozhukov, Iván Werning, and Michael D. Whinston (2020). “A
2696 multi-risk SIR model with optimally targeted lockdown”. In: *NBER Working Paper Series*.
- 2697 Adam, Lukáš, Rostislav Horčík, Tomáš Kasl, and Tomáš Kroupa (2021). “Double oracle algorithm
2698 for computing equilibria in continuous games”. In: *AAAI Conference on Artificial Intelligence*.
- 2699 Adamo, Tim and Alexander Matros (2009). “A Blotto game with incomplete information”. In:
2700 *Economics Letters*.
- 2701 Aggarwal, Gagan and Jason D. Hartline (2006). “Knapsack auctions”. In: *SODA*.
- 2702 Agrawal, Akshay, Robin Verschueren, Steven Diamond, and Stephen Boyd (2018). “A rewriting
2703 system for convex optimization problems”. In: *Journal of Control and Decision*.
- 2704 Alfriend, Kyle T., Srinivas R. Vadali, Pini Gurfil, Jonathan P. How, and Louis S. Breger (2009).
2705 *Spacecraft formation flying: dynamics, control and navigation*. Elsevier.
- 2706 Arnott, Richard, André De Palma, and Robin Lindsey (1990). “Economics of a bottleneck”. In:
2707 *Journal of urban economics*.
- 2708 Aumann, Robert (1964). “Markets with a continuum of traders”. In: *Econometrica*.
- 2709 — (1974). “Subjectivity and correlation in randomized strategies”. In: *Journal of Mathematical*
2710 *Economics*.
- 2711 Bachmann, Paul (1894). *Die Analytische Zahlentheorie (Analytic Number Theory)*. Teubner.
- 2712 Bäck, Thomas (1996). *Evolutionary algorithms in theory and practice*. Oxford University Press.
- 2713 Bäck, Thomas, David B. Fogel, and Zbigniew Michalewicz (1997). *Handbook of evolutionary compu-*
2714 *tation*. IOP Publishing Ltd.
- 2715 Bae, Juhan and Roger B. Grosse (2020). “Delta-STN: efficient bilevel optimization for neural networks
2716 using structured response Jacobians”. In: *Conference on Neural Information Processing Systems*
2717 *(NeurIPS)*.
- 2718 Bailey, Bolton and Matus Telgarsky (2018). “Size-noise tradeoffs in generative networks”. In: *Confer-*
2719 *ence on Neural Information Processing Systems (NeurIPS)*.
- 2720 Balduzzi, David, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore
2721 Graepel (2018). “The mechanics of n-player differentiable games”. In: *International Conference*
2722 *on Machine Learning (ICML)*.
- 2723 Bao, Xuchan and Guodong Zhang (2022). “Finding and only finding local Nash equilibria by both
2724 pretending to be a follower”. In: *Workshop on Gamification and Multiagent Solutions*.
- 2725 Bardi, Martino and Italo Capuzzo-Dolcetta (1997). *Optimal control and viscosity solutions of*
2726 *Hamilton-Jacobi-Bellman Equations*. Birkhäuser Boston.
- 2727 Başar, Tamer and Geert Jan Olsder (1999). *Dynamic noncooperative game theory*. SIAM.
- 2728 Baye, Michael R., Dan Kovenock, and Casper G. de Vries (1996). “The all-pay auction with complete
2729 information”. In: *Economic Theory*.

- 2730 Berahas, Albert, Liyuan Cao, Krzysztof Choromanski, and Katya Scheinberg (2022). “A theoret-
2731 ical and empirical comparison of gradient approximations in derivative-free optimization”. In:
2732 *Foundations of Computational Mathematics*.
- 2733 Berg, Jordan, Amy Greenwald, Victor Naroditskiy, and Eric Sodomka (2010). “A knapsack-based
2734 approach to bidding in ad auctions”. In: *European Conference on Artificial Intelligence (ECAI)*.
2735 IOS Press.
- 2736 Berger, Ulrich (2007). “Brown’s original fictitious play”. In: *Journal of Economic Theory (JET)*.
- 2737 Bergman, L. M. and I. N. Fokin (1998). “On separable non-cooperative zero-sum games”. In:
2738 *Optimization*.
- 2739 Berner, Christopher, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębniak, Christy
2740 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Jozefowicz, Scott Gray,
2741 Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d.O. Pinto, Jonathan Raiman,
2742 Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip
2743 Wolski, and Susan Zhang (2019). “Dota 2 with large scale deep reinforcement learning”. In:
2744 *arXiv:1912.06680*.
- 2745 Berridge, Steffan and Jacek B. Krawczyk (1997). *Relaxation algorithms in finding Nash equilibria*.
2746 Tech. rep. Victoria University of Wellington.
- 2747 Berry, Andrew C. (1941). “The accuracy of the Gaussian approximation to the sum of independent
2748 variates”. In: *Transactions of the American Mathematical Society*.
- 2749 Bewley, Truman F. (1972). “Existence of equilibria in economies with infinitely many commodities”.
2750 In: *Journal of Economic Theory (JET)*.
- 2751 Bichler, Martin, Max Fichtl, and Matthias Oberlechner (2023). “Computing Bayes–Nash equilibrium
2752 strategies in auction games via simultaneous online dual averaging”. In: *Operations Research*.
- 2753 Bichler, Martin, Maximilian Fichtl, Stefan Heidekrüger, Nils Kohring, and Paul Sutterer (2021).
2754 “Learning equilibria in symmetric auction games using artificial neural networks”. In: *Nature*
2755 *Machine Intelligence*.
- 2756 Bichler, Martin, Nils Kohring, and Stefan Heidekrüger (2023). “Learning equilibria in asymmetric
2757 auction games”. In: *INFORMS Journal on Computing*.
- 2758 Bochner, Salomon (1933). “Integration von Funktionen, deren Werte die Elemente eines Vektorraumes
2759 sind”. In: *Fundamenta Mathematicae*.
- 2760 Boix-Adserà, Enric, Benjamin L. Edelman, and Siddhartha Jayanti (2021). “The multiplayer Colonel
2761 Blotto game”. In: *Games and Economic Behavior (GEB)*.
- 2762 Bolton, Gary E. and Axel Ockenfels (2000). “ERC: A theory of equity, reciprocity, and competition”.
2763 In: *American Economic Review (AER)*.
- 2764 Bonabeau, Eric (2002). “Agent-based modeling: methods and techniques for simulating human
2765 systems”. In: *Proceedings of the National Academy of Sciences (PNAS)*.
- 2766 Borel, Émile (1921). “The theory of play and integral equations with skew symmetric kernels”. In:
2767 *Comptes Rendus de l’Académie des Sciences*.
- 2768 Borel, Émile and Jean Ville (1938). *Traité du calcul des probabilités et ses applications*. Gauthier-
2769 Villars.
- 2770 Borenstein, Severin, James Bushnell, and Christopher R. Knittel (1999). “Market power in electricity
2771 markets: beyond concentration measures”. In: *The Energy Journal*.
- 2772 Bradbury, James, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
2773 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
2774 Zhang (2018). *JAX: composable transformations of Python+NumPy programs*.

- 2775 Brandt, Felix, Tuomas Sandholm, and Yoav Shoham (2007). “Spiteful bidding in sealed-bid auctions”.
2776 In: *International Joint Conference on Artificial Intelligence (IJCAI)*.
- 2777 Bravo, Mario, David Leslie, and Panayotis Mertikopoulos (2018). “Bandit learning in concave
2778 N-person games”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 2779 Bressan, Alberto (2015). “Conservation law models for traffic flow on a network of roads”. In:
2780 *Networks and Heterogeneous Media*.
- 2781 Brouwer, Luitzen Egbertus Jan (1911). “Über abbildung von mannigfaltigkeiten”. In: *Mathematische
2782 annalen*.
- 2783 Brown, George W. (1951). “Iterative solution of games by fictitious play”. In: *Activity Analysis of
2784 Production and Allocation*. Wiley.
- 2785 Brown, Noam and Tuomas Sandholm (2018). “Superhuman AI for heads-up no-limit poker: Libratus
2786 beats top professionals”. In: *Science*.
- 2787 — (2019). “Superhuman AI for multiplayer poker”. In: *Science*.
- 2788 Bunyakovsky, Viktor Y. (1859). “Sur quelques inégalités concernant les intégrales aux différences
2789 finis”. In: *Mémoires de l’Académie Impériale des Sciences de St. Pétersbourg*.
- 2790 Busoniu, Lucian, Robert Babuska, and Bart De Schutter (2008). “A comprehensive survey of
2791 multiagent reinforcement learning”. In: *IEEE Transactions on Systems, Man, and Cybernetics,
2792 Part C (Applications and Reviews)*.
- 2793 Cai, Han, Kan Ren, Weinan Zhang, Kleantlis Malialis, Jun Wang, Yong Yu, and Defeng Guo (2017).
2794 “Real-time bidding by reinforcement learning in display advertising”. In: *ACM International
2795 Conference on Web Search and Data Mining (WSDM)*.
- 2796 Cai, Yang, Ozan Candogan, Constantinos Daskalakis, and Christos Papadimitriou (2016). “Zero-sum
2797 polymatrix games: a generalization of minmax”. In: *Mathematics of Operations Research*.
- 2798 Cai, Yang and Constantinos Daskalakis (2011). “On minmax theorems for multiplayer games”. In:
2799 *ACM-SIAM Symposium on Discrete Algorithms (SODA)*.
- 2800 Caines, Peter and Minyi Huang (2021). “Graphon mean field games and their equations”. In: *SIAM
2801 Journal on Control and Optimization*.
- 2802 Caines, Peter, Minyi Huang, and Roland Malhamé (2018). “Mean field games”. In: *Handbook of
2803 Dynamic Game Theory*. Springer International Publishing.
- 2804 Cauchy, Augustin Louis Baron (1821a). *Cours d’analyse de l’École Royale Polytechnique*. Imprimerie
2805 royale.
- 2806 Cauchy, Augustin-Louis (1821b). “Sur les formules qui résultent de l’emploi du signe et sur $>$ ou $<$, et
2807 sur les moyennes entre plusieurs quantités”. In: *Cours d’Analyse, 1er Partie: Analyse algébrique*.
- 2808 Chan, Patrick and Ronnie Sircar (2015). “Bertrand and Cournot mean field games”. In: *Applied
2809 Mathematics & Optimization*.
- 2810 Charness, Gary and Matthew Rabin (2002). “Understanding social preferences with simple tests”.
2811 In: *The Quarterly Journal of Economics*.
- 2812 Chen, Bill and Jerrod Ankenman (2006). *The Mathematics of Poker*. ConJelCo.
- 2813 Clevert, Djork-Arné, Thomas Unterthiner, and Sepp Hochreiter (2016). “Fast and accurate deep
2814 network learning by exponential linear units (ELUs)”. In: *International Conference on Learning
2815 Representations (ICLR)*.
- 2816 Cloud, Alex, Albert Wang, and Wesley Kerr (2023). “Anticipatory fictitious play”. In: *International
2817 Joint Conference on Artificial Intelligence (IJCAI)*.
- 2818 Cornes, Richard and Roger Hartley (2007). “Aggregative public good games”. In: *Journal of Public
2819 Economic Theory*.

- 2820 Coulom, Rémi (2007). “Efficient selectivity and backup operators in Monte–Carlo tree search”. In:
2821 *Computers and Games*. Springer.
- 2822 Cournot, Antoine Augustin (1838). *Recherches sur les principes mathématiques de la théorie des*
2823 *richesses*. Hachette.
- 2824 — (1863). *Principes de la théorie des richesses*. Hachette.
- 2825 Crouse, David F. (2016). “On implementing 2D rectangular assignment algorithms”. In: *IEEE*
2826 *Transactions on Aerospace and Electronic Systems*.
- 2827 Cui, Kai and Heinz Koepl (2022). “Learning graphon mean field games and approximate Nash
2828 equilibria”. In: *International Conference on Learning Representations (ICLR)*.
- 2829 Cuturi, Marco, Laetitia Meng-Papaxanthos, Yingtao Tian, Charlotte Bunne, Geoff Davis, and Olivier
2830 Teboul (2022). “Optimal transport tools (OTT): a JAX Toolbox for all things Wasserstein”. In:
2831 *arXiv:2201.12324*.
- 2832 Cybenko, George (1989). “Approximation by superpositions of a sigmoidal function”. In: *Mathematics*
2833 *of control, signals and systems*.
- 2834 Dasgupta, Partha and Eric Maskin (1986). “The existence of equilibrium in discontinuous economic
2835 games 1: theory”. In: *Review of Economic Studies*.
- 2836 Daskalakis, Constantinos, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng (2018). “Training
2837 GANs with optimism”. In: *International Conference on Learning Representations (ICLR)*.
- 2838 DeepMind, Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter
2839 Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Claudio
2840 Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo
2841 Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena
2842 Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John Quan,
2843 Roman Ring, Francisco Ruiz, Alvaro Sanchez, Laurent Sartran, Rosalia Schneider, Eren Sezener,
2844 Stephen Spencer, Srivatsan Srinivasan, Miloš Stanojević, Wojciech Stokowiec, Luyu Wang,
2845 Guangyao Zhou, and Fabio Viola (2020). *The DeepMind JAX ecosystem*.
- 2846 Deng, Li (2012). “The MNIST database of handwritten digit images for machine learning research”.
2847 In: *IEEE Signal Processing Magazine*.
- 2848 Diamond, Steven and Stephen Boyd (2016). “CVXPY: a Python-embedded modeling language for
2849 convex optimization”. In: *Journal of Machine Learning Research*.
- 2850 Dockner, Engelbert, Steffen Jorgensen, Van Long Ngo, and Gerhard Sorger (2000). *Differential*
2851 *games in economics and management science*. Cambridge University Press.
- 2852 Domingo-Enrich, Carles (2019). “Games in machine learning: differentiable n-player games and
2853 structured planning”. MA thesis. Universitat Politècnica de Catalunya.
- 2854 Dormand, John R. and Peter J. Prince (1980). “A family of embedded Runge–Kutta formulae”. In:
2855 *Journal of computational and applied mathematics*.
- 2856 Dror, Moshe (1989). “Simple proof for Goofspiel: the game of pure strategy”. In: *Advances in Applied*
2857 *Probability*.
- 2858 Duane, Simon, Anthony D. Kennedy, Brian J. Pendleton, and Duncan Roweth (1987). “Hybrid
2859 Monte Carlo”. In: *Physics letters B*.
- 2860 Duchi, John C., Peter L. Bartlett, and Martin J. Wainwright (2012). “Randomized smoothing for
2861 stochastic optimization”. In: *SIAM Journal on Optimization*.
- 2862 Duchi, John C., Michael I. Jordan, Martin J. Wainwright, and Andre Wibisono (2015). “Optimal
2863 rates for zero-order convex optimization: the power of two function evaluations”. In: *IEEE*
2864 *Transactions on Information Theory*.

- 2865 Dunford, Nelson (1935). “Integration in general analysis”. In: *Transactions of the American Mathe-*
2866 *matical Society*.
- 2867 Dütting, Paul, Vasilis Gkatzelis, and Tim Roughgarden (2014). “The performance of deferred-
2868 acceptance auctions”. In: *ACM Conference on Economics and Computation (EC)*.
- 2869 Edelman, Benjamin, Michael Ostrovsky, and Michael Schwarz (2007). “Internet advertising and the
2870 generalized second-price auction: selling billions of dollars worth of keywords”. In: *American*
2871 *Economic Review (AER)*.
- 2872 Efron, Bradley (1979). “Bootstrap methods: another look at the jackknife”. In: *The Annals of*
2873 *Statistics*.
- 2874 — (1987). “Better bootstrap confidence intervals”. In: *Journal of the American Statistical Association*
2875 *(JASA)*.
- 2876 Eiben, Agoston E. and James E. Smith (2003). *Introduction to evolutionary computing*. Springer.
- 2877 Eisenberg, Larry and Thomas H. Noe (2001). “Systemic risk in financial systems”. In: *Management*
2878 *Science*.
- 2879 Elliott, Matthew, Benjamin Golub, and Matthew O. Jackson (2014). “Financial networks and
2880 contagion”. In: *American Economic Review (AER)*.
- 2881 Epstein, Joshua and Robert Axtell (1996). *Growing artificial societies: social science from the bottom*
2882 *up*. Brookings Institution Press.
- 2883 Fehr, Ernst and Klaus M. Schmidt (1999). “A theory of fairness, competition, and cooperation”. In:
2884 *The Quarterly Journal of Economics*.
- 2885 Feldman, William, Inwon Kim, and Aaron Zeff Palmer (2024). “The sharp interface limit of an Ising
2886 game”. In: *ESAIM: Control, Optimisation and Calculus of Variations*.
- 2887 Feng, Ruili, Deli Zhao, and Zheng-Jun Zha (2021). “Understanding noise injection in GANs”. In:
2888 *International Conference on Machine Learning (ICML)*.
- 2889 Ferguson, Thomas S. and Costis Melolidakis (2001). “Games with finite resources”. In: *International*
2890 *Journal of Game Theory (IJGT)*.
- 2891 Fiez, Tanner, Benjamin Chasnov, and Lillian Ratliff (2020). “Implicit learning dynamics in stackelberg
2892 games: equilibria characterization, convergence analysis, and empirical study”. In: *International*
2893 *Conference on Machine Learning (ICML)*.
- 2894 Fiez, Tanner, Benjamin Chasnov, and Lillian J. Ratliff (2019). “Convergence of learning dynamics in
2895 Stackelberg games”. In: *arXiv:1906.01217*.
- 2896 Fiez, Tanner, Chi Jin, Praneeth Netrapalli, and Lillian J. Ratliff (2022). “Minimax optimization
2897 with smooth algorithmic adversaries”. In: *International Conference on Learning Representations*
2898 *(ICLR)*.
- 2899 Flåm, Sjur Didrik and Anatoly S. Antipin (1996). “Equilibrium programming using proximal-like
2900 algorithms”. In: *Mathematical Programming*.
- 2901 Flåm, Sjur Didrik and Andrzej Ruszczyński (2008). “Finding normalized equilibrium in convex-
2902 concave games”. In: *International Game Theory Review*.
- 2903 Foerster, Jakob N., Richard Y. Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and
2904 Igor Mordatch (2018). “Learning with opponent-learning awareness”. In: *International Conference*
2905 *on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- 2906 French, Robert M. (1999). “Catastrophic forgetting in connectionist networks”. In: *Trends in cognitive*
2907 *sciences*.
- 2908 Friedman, Avner (1971a). “Computation of saddle points for differential games of pursuit and evasion”.
2909 In: *Archive for Rational Mechanics and Analysis*.
- 2910 — (1971b). *Differential games*. Wiley-Interscience.

- 2911 Fu, Michael C., ed. (2015). *Handbook of simulation optimization*. Springer.
- 2912 Fudenberg, Drew and Jean Tirole (1991). *Game theory*. MIT Press.
- 2913 Gadetsky, Artyom, Kirill Struminsky, Christopher Robinson, Novi Quadrianto, and Dmitry Vetrov
2914 (2020). “Low-variance black-box gradient estimates for the Plackett–Luce distribution”. In: *AAAI
2915 Conference on Artificial Intelligence*.
- 2916 Gal, Yarín and Zoubin Ghahramani (2016). “Dropout as a Bayesian approximation: representing
2917 model uncertainty in deep learning”. In: *International Conference on Machine Learning (ICML)*.
- 2918 Galam, Serge and Bernard Walliser (2010). “Ising model versus normal form game”. In: *Physica A*.
- 2919 Ganzfried, Sam (2021). “Algorithm for computing approximate Nash equilibrium in continuous
2920 games with application to continuous Blotto”. In: *Games*.
- 2921 Ganzfried, Sam and Tuomas Sandholm (2010). “Computing equilibria by incorporating qualita-
2922 tive models”. In: *International Conference on Autonomous Agents and Multi-Agent Systems
2923 (AAMAS)*.
- 2924 Gemp, Ian and Sridhar Mahadevan (2018). “Global convergence to the equilibrium of GANs using
2925 variational inequalities”. In: *arXiv:1808.01531*.
- 2926 Gemp, Ian, Rahul Savani, Marc Lanctot, Yoram Bachrach, Thomas Anthony, Richard Everett,
2927 Andrea Tacchetti, Tom Eccles, and János Kramár (2022). “Sample-based approximation of Nash
2928 in large many-player games via gradient descent”. In: *International Conference on Autonomous
2929 Agents and Multi-Agent Systems (AAMAS)*.
- 2930 Ghosh, Papiya and Rajendra P. Kundu (2019). “Best-shot network games with continuous action
2931 space”. In: *Research in Economics*.
- 2932 Gilles, Christian and Stephen F. LeRoy (1992). “Bubbles and charges”. In: *International Economic
2933 Review*.
- 2934 Glicksberg, Irving Leonard (1952). “A further generalization of the Kakutani fixed point theorem,
2935 with application to Nash equilibrium points”. In: *Proceedings of the American Mathematical
2936 Society*.
- 2937 Glicksberg, Irving Leonard and Oliver Alfred Gross (1953). “Notes on games over the square”. In:
2938 *Contributions to the Theory of Games*.
- 2939 Glorot, Xavier and Yoshua Bengio (2010). “Understanding the difficulty of training deep feedforward
2940 neural networks”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- 2941 Goktas, Denizalp and Amy Greenwald (2021). “Convex-concave min-max Stackelberg games”. In:
2942 *Conference on Neural Information Processing Systems (NeurIPS)*.
- 2943 — (2022a). “Exploitability minimization in games and beyond”. In: *Conference on Neural Information
2944 Processing Systems (NeurIPS)*.
- 2945 — (2022b). “Gradient descent ascent in min-max Stackelberg games”. In: *arXiv:2208.09690*.
- 2946 Goktas, Denizalp, Jiayi Zhao, and Amy Greenwald (2022a). “Robust no-regret learning in min-
2947 max Stackelberg games”. In: *International Conference on Autonomous Agents and Multi-Agent
2948 Systems (AAMAS)*.
- 2949 Goktas, Denizalp, Sadie Zhao, and Amy Greenwald (2022b). “Zero-sum stochastic Stackelberg
2950 games”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 2951 Goldblum, Micah, Marc Finzi, Keefer Rowan, and Andrew Gordon Wilson (2024). “Position: the no
2952 free lunch theorem, Kolmogorov complexity, and the role of inductive biases in machine learning”.
2953 In: *International Conference on Machine Learning (ICML)*.
- 2954 Golowich, Noah, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar (2020). “Last
2955 iterate is slower than averaged iterate in smooth convex-concave saddle point problems”. In:
2956 *Conference on Learning Theory (COLT)*.

- 2957 Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
2958 Aaron Courville, and Yoshua Bengio (2014). “Generative adversarial nets”. In: *Conference on*
2959 *Neural Information Processing Systems (NeurIPS)*.
- 2960 Grimes, Mark and Moshe Dror (2013). “Observations on strategies for Goofspiel”. In: *IEEE CIG*.
- 2961 Grimm, Volker, Eloy Revilla, Uta Berger, Florian Jeltsch, Wolf M. Mooij, Steven F. Railsback,
2962 Hans-Hermann Thulke, Jacob Weiner, Thorsten Wiegand, and Donald L. DeAngelis (2005).
2963 “Pattern-oriented modeling of agent-based complex systems: lessons from ecology”. In: *Science*.
- 2964 Grnarova, Paulina, Yannic Kilcher, Kfir Y. Levy, Aurelien Lucchi, and Thomas Hofmann (2021).
2965 “Generative minimization networks: training GANs without competition”. In: *arXiv:2103.12685*.
- 2966 Grnarova, Paulina, Kfir Y. Levy, Aurelien Lucchi, Nathanael Perraudin, Ian Goodfellow, Thomas
2967 Hofmann, and Andreas Krause (2019). “A domain agnostic measure for monitoring and evaluating
2968 GANs”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 2969 Gross, Oliver Alfred and R. A. Wagner (1950). *A continuous Colonel Blotto game*. RAND Corpora-
2970 tion.
- 2971 Grover, Aditya, Eric Wang, Aaron Zweig, and Stefano Ermon (2019). “Stochastic optimization of
2972 sorting networks via continuous relaxations”. In: *International Conference on Learning Representations*.
2973
- 2974 Guo, Xin, Anran Hu, Renyuan Xu, and Junzi Zhang (2019). “Learning mean-field games”. In:
2975 *Conference on Neural Information Processing Systems (NeurIPS)*.
- 2976 Gürkan, Gül and Jong-Shi Pang (2009). “Approximations of Nash equilibria”. In: *Mathematical*
2977 *Programming*.
- 2978 Ha, David, Andrew M. Dai, and Quoc V. Le (2017). “HyperNetworks”. In: *International Conference*
2979 *on Learning Representations (ICLR)*.
- 2980 Hairer, Ernst, Syvert P. Nørsett, and Gerhard Wanner (1993). *Solving ordinary differential equations*
2981 *I: nonstiff problems*. Springer-Verlag.
- 2982 Hall, J., I. Galabova, L. Gottwald, and M. Feldmeier (2023). *HiGHS—high performance software for*
2983 *linear optimization*.
- 2984 Hamilton, William D. (1964). “The genetical evolution of social behaviour. I and II”. In: *Journal of*
2985 *Theoretical Biology*.
- 2986 Harris, Charles R., K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen,
2987 David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert
2988 Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane,
2989 Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard,
2990 Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant
2991 (2020). “Array programming with NumPy”. In: *Nature*.
- 2992 Harsanyi, John (1967). “Games with incomplete information played by Bayesian players”. In:
2993 *Management Science*.
- 2994 Hastings, Wilfred (1970). “Monte Carlo sampling methods using Markov chains and their applications”.
2995 In: *Biometrika*.
- 2996 He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (2015). “Delving deep into rectifiers:
2997 surpassing human-level performance on ImageNet classification”. In: *International Conference on*
2998 *Computer Vision (ICCV)*.
- 2999 Heek, Jonathan, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas
3000 Steiner, and Marc van Zee (2023). *Flax: a neural network library and ecosystem for JAX*.
- 3001 Heinrich, Johannes, Marc Lanctot, and David Silver (2015). “Fictitious self-play in extensive-form
3002 games”. In: *International Conference on Machine Learning (ICML)*.

- 3003 Heinrich, Johannes and David Silver (2016). “Deep reinforcement learning from self-play in imperfect-
3004 information games”. In: *arXiv:1603.01121*.
- 3005 Heusinger, Anna Von and Christian Kanzow (2009a). “Optimization reformulations of the generalized
3006 Nash equilibrium problem using Nikaido–Isoda-type functions”. In: *Computational Optimization
3007 and Applications*.
- 3008 — (2009b). “Relaxation methods for generalized Nash equilibrium problems with inexact line search”.
3009 In: *Journal of Optimization Theory and Applications (JOTA)*.
- 3010 Hoffman, Matthew D. and Andrew Gelman (2014). “The No-U-Turn sampler: adaptively setting
3011 path lengths in Hamiltonian Monte Carlo”. In: *Journal of Machine Learning Research (JMLR)*.
- 3012 Hornik, Kurt (1991). “Approximation capabilities of multilayer feedforward networks”. In: *Neural
3013 networks*.
- 3014 Hou, Jian, Zong-Chuan Wen, and Qing Chang (2018). “An unconstrained optimization reformulation
3015 for the Nash game”. In: *Journal of Interdisciplinary Mathematics (JIM)*.
- 3016 Hsieh, Yu-Guan, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos (2019). “On the con-
3017 vergence of single-call stochastic extra-gradient methods”. In: *Conference on Neural Information
3018 Processing Systems (NeurIPS)*.
- 3019 Hsieh, Ya-Ping, Panayotis Mertikopoulos, and Volkan Cevher (2021). “The limits of min-max
3020 optimization algorithms: convergence to spurious non-critical sets”. In: *International Conference
3021 on Machine Learning (ICML)*.
- 3022 Hu, Haigen, Xiaoyuan Wang, Yan Zhang, Qi Chen, and Qiu Guan (2024). “A comprehensive survey
3023 on contrastive learning”. In: *Neurocomputing*.
- 3024 Hu, Yaohua, Xiaoqi Yang, and Chee-Khian Sim (2015). “Inexact subgradient methods for quasi-convex
3025 optimization problems”. In: *European Journal of Operational Research*.
- 3026 Huang, Chin-Wei, David Krueger, Alexandre Lacoste, and Aaron Courville (2018). “Neural autore-
3027 gressive flows”. In: *International Conference on Machine Learning (ICML)*.
- 3028 Huang, Kevin XD and Jan Werner (2000). “Asset price bubbles in Arrow–Debreu and sequential
3029 equilibrium”. In: *Economic Theory*.
- 3030 Huang, Minyi, Peter E. Caines, and Roland P. Malhamé (2007). “Large-population cost-coupled
3031 LQG problems with nonuniform agents: individual-mass behavior and decentralized epsilon-Nash
3032 equilibria”. In: *IEEE Transactions on Automatic Control*.
- 3033 Huang, Minyi, Roland P. Malhamé, and Peter E. Caines (2006). “Large population stochastic dynamic
3034 games: closed-loop McKean–Vlasov systems and the Nash certainty equivalence principle”. In:
3035 *Communications in Information and Systems*.
- 3036 Huangfu, Qi and J. A. Julian Hall (2018). “Parallelizing the dual revised simplex method”. In:
3037 *Mathematical Programming Computation*.
- 3038 Hunter, John (2007). “Matplotlib: a 2D graphics environment”. In: *Computing in Science & Engi-
3039 neering*.
- 3040 Isaacs, Rufus (1965). *Differential games: a mathematical theory with applications to warfare and
3041 pursuit, control and optimization*. John Wiley and Sons.
- 3042 Ising, Ernst (1925). “Beitrag zur theorie des ferromagnetismus”. In: *Zeitschrift für Physik*.
- 3043 Itô, Kiyosi (1944). “Stochastic integral”. In: *Proceedings of the Imperial Academy*.
- 3044 Jaiswal, Ashish, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia
3045 Makedon (2020). “A survey on contrastive self-supervised learning”. In: *Technologies*.
- 3046 Ji, Kaiyi, Junjie Yang, and Yingbin Liang (2021). “Bilevel optimization: Convergence analysis and
3047 enhanced design”. In: *International conference on machine learning*. Proceedings of Machine
3048 Learning Research (PMLR).

- 3049 Jin, Chi, Praneeth Netrapalli, and Michael I. Jordan (2020). “What is local optimality in nonconvex-
3050 nonconcave minimax optimization?” In: *International Conference on Machine Learning (ICML)*.
- 3051 Jonker, Roy and Ton Volgenant (1988). “A shortest augmenting path algorithm for dense and sparse
3052 linear assignment problems”. In: *DGOR/NSOR*.
- 3053 Jordan, James S. (1993). “Three problems in learning mixed-strategy Nash equilibria”. In: *Games
3054 and Economic Behavior (GEB)*.
- 3055 Kagel, John H. and Dan Levin (1993). “Independent private value auctions: bidder behaviour in
3056 first-, second- and third-price auctions with varying numbers of bidders”. In: *The Economic
3057 Journal*.
- 3058 Kakutani, Shizuo (1941). “A generalization of Brouwer’s fixed point theorem”. In: *Duke Mathematical
3059 Journal*.
- 3060 Kamra, Nitin, Fei Fang, Debarun Kar, Yan Liu, and Milind Tambe (2017). “Handling continuous
3061 space security games with neural networks”. In: *International Joint Conference on Artificial
3062 Intelligence (IJCAI)*.
- 3063 Kamra, Nitin, Umang Gupta, Fei Fang, Yan Liu, and Milind Tambe (2018). “Policy learning for
3064 continuous space security games using neural networks”. In: *AAAI Conference on Artificial
3065 Intelligence*.
- 3066 Kamra, Nitin, Umang Gupta, Kai Wang, Fei Fang, Yan Liu, and Milind Tambe (2019). “DeepFP
3067 for finding Nash equilibrium in continuous action spaces”. In: *Decision and Game Theory for
3068 Security*.
- 3069 Kar, Debarun, Thanh H. Nguyen, Fei Fang, Matthew Brown, Arunesh Sinha, Milind Tambe, and
3070 Albert Xin Jiang (2017). “Trends and applications in Stackelberg security games”. In: *Handbook
3071 of dynamic game theory*. Springer International Publishing.
- 3072 Kessler, Donald J. and Burton G. Cour-Palais (1978). “Collision frequency of artificial satellites: the
3073 creation of a debris belt”. In: *Journal of Geophysical Research: Space Physics*.
- 3074 Le-Khac, Phuc, Graham Healy, and Alan Smeaton (2020). “Contrastive representation learning: a
3075 framework and review”. In: *IEEE Access*.
- 3076 Khan, Ali (1985). “Equilibrium points of nonatomic games over a nonreflexive Banach space”. In:
3077 *Journal of Approximation Theory*.
- 3078 — (1986). “Equilibrium points of nonatomic games over a Banach space”. In: *Transactions of the
3079 American Mathematical Society*.
- 3080 Khan, Ali and Nikolaos Papageorgiou (1987a). “On Cournot–Nash equilibria in generalized qualitative
3081 games with a continuum of players”. In: *Nonlinear Analysis: Theory, Methods & Applications*.
- 3082 — (1987b). “On Cournot–Nash equilibria in generalized qualitative games with an atomless measure
3083 space of agents”. In: *Proceedings of the American Mathematical Society*.
- 3084 Khan, Ali and Yeneng Sun (2002). “Non-cooperative games with many players”. In: *Handbook of
3085 Game Theory with Economic Applications*.
- 3086 Kidger, Patrick (2021). “On neural differential equations”. PhD thesis. University of Oxford.
- 3087 Kidger, Patrick and Terry Lyons (2020). “Universal approximation with deep narrow networks”. In:
3088 *Conference on Learning Theory (COLT)*.
- 3089 Kim, Taesung, Karel Prikry, and Nicholas Yannellis (1989). “Equilibria in abstract economies with
3090 a measure space of agents and with an infinite dimensional strategy space”. In: *Journal of
3091 Approximation Theory*.
- 3092 Kirkpatrick, James, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A.
3093 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis,

- 3094 Claudia Clopath, Dharshan Kumaran, and Raia Hadsell (2017). “Overcoming catastrophic
3095 forgetting in neural networks”. In: *Proceedings of the National Academy of Sciences (PNAS)*.
- 3096 Kiwiel, Krzysztof C. (2004). “Convergence of approximate and incremental subgradient methods for
3097 convex optimization”. In: *SIAM Journal on Optimization*.
- 3098 Kizilkale, A. Caağatay, Rabih Salhab, and Roland P. Malhamé (2019). “An integral control formula-
3099 tion of mean field game based large scale coordination of loads in smart grids”. In: *Automatica*.
- 3100 Kleinert, Thomas, Martine Labbé, Ivana Ljubić, and Martin Schmidt (2021). “A survey on mixed-
3101 integer programming techniques in bilevel optimization”. In: *EURO Journal on Computational
3102 Optimization*.
- 3103 Knuth, Donald E. (1976). “Big omicron and big omega and big theta”. In: *ACM Sigact News*.
- 3104 Koranne, Sandeep (2010). “Hierarchical data format 5: HDF5”. In: *Handbook of open source tools*.
3105 Springer.
- 3106 Korpelevich, Galina M. (1976). “The extragradient method for finding saddle points and other
3107 problems”. In: *Ekonomika i Matematicheskie Metody*.
- 3108 Korzhyk, Dmytro, Zhengyu Yin, Christopher Kiekintveld, Vincent Conitzer, and Milind Tambe
3109 (2011). “Stackelberg vs. Nash in security games”. In: *Journal of Artificial Intelligence Research
3110 (JAIR)*.
- 3111 Kovenock, Dan and Brian Roberson (2011). “A Blotto game with multi-dimensional incomplete
3112 information”. In: *Economics Letters*.
- 3113 — (2021). “Generalizations of the General Lotto and Colonel Blotto games”. In: *Economic Theory*.
- 3114 Krawczyk, Jacek B. (2005). “Coupled constraint Nash equilibria in environmental games”. In: *Resource
3115 and Energy Economics*.
- 3116 Krawczyk, Jacek B. and Stanislav Uryasev (2000). “Relaxation algorithms to find Nash equilibria
3117 with economic applications”. In: *Environmental Modeling & Assessment*.
- 3118 Krishna, Vijay (2002). *Auction theory*. Academic Press.
- 3119 Kroer, Christian and Tuomas Sandholm (2015). “Discretization of continuous action spaces in
3120 extensive-form games”. In: *International Conference on Autonomous Agents and Multi-Agent
3121 Systems (AAMAS)*.
- 3122 Kroupa, Tomáš and Tomáš Votroubek (2023). “Multiple oracle algorithm to solve continuous games”.
3123 In: *Decision and Game Theory for Security*.
- 3124 Kuhn, Harold W. (1955). “The Hungarian method for the assignment problem”. In: *Naval research
3125 logistics quarterly*.
- 3126 Kuhn, Harold William (1950). “A simplified two-person poker”. In: *Contributions to the Theory of
3127 Games*. Princeton University Press.
- 3128 Kuipers, Lauwerens and Harald Niederreiter (1974). *Uniform distribution of sequences*. Wiley-
3129 Interscience.
- 3130 Kurzweil, Jaroslav (1957). “Generalized ordinary differential equations and continuous dependence
3131 on a parameter”. In: *Czechoslovak Mathematical Journal*.
- 3132 Lanctot, Marc, Viliam Lisý, and Mark H. M. Winands (2014). “Monte Carlo tree search in simulta-
3133 neous move games with applications to Goofspiel”. In: *Computer Games*.
- 3134 Lanctot, Marc, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Jason Upchurch, Julien
3135 Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, Daniel Hennes,
3136 Dustin Morrill, Paul Muller, Timo Ewalds, Ryan Sutherland, Chris Anderson, Ulrich Siegler,
3137 Sandra Bechtle, Sarah McAleer, Richard Dehouck, David Kasenberg, Michael Ostercamp, Jordan
3138 Sahli, Victor Shkurti, Venugopal Lakshminarayanan, Panagiotis Kougiouris, Justin Fu, Houghton

- 3139 Bell, Abhishek Sen, and Stephen Knight (2019). “OpenSpiel: a framework for reinforcement
3140 learning in games”. In: *arXiv:1908.09453*.
- 3141 Lanctot, Marc, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat,
3142 David Silver, and Thore Graepel (2017). “A unified game-theoretic approach to multiagent
3143 reinforcement learning”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 3144 Landau, Edmund (1909). *Handbuch der Lehre von der Verteilung der Primzahlen (Handbook of the
3145 Theory of the Distribution of Prime Numbers)*. BG Teubner.
- 3146 Lasry, Jean-Michel and Pierre-Louis Lions (2007). “Mean field games”. In: *Japanese Journal of
3147 Mathematics*.
- 3148 Laurière, Mathieu, Sarah Perrin, Julien Pérolat, Sertan Girgin, Paul Muller, Romuald Élie, Matthieu
3149 Geist, and Olivier Pietquin (2022). “Learning mean field games: a survey”. In: *arXiv:2205.12944*.
- 3150 Lebesgue, Henri (1902). “Intégrale, longueur, aire”. In: *Annali di Matematica Pura ed Applicata
3151 (1898-1922)*.
- 3152 LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton (2015). “Deep learning”. In: *Nature*.
- 3153 Lenc, Karel, Erich Elsen, Tom Schaul, and Karen Simonyan (2019). “Non-differentiable supervised
3154 learning with evolution strategies and hybrid methods”. In: *arXiv:1906.03139*.
- 3155 Lenz, Wilhelm (1920). “Beitrag zum Verständnis der magnetischen Erscheinungen in festen Körpern”.
3156 In: *Physikalische Zeitschrift*.
- 3157 Leonidov, Andrey, Alexey Savvateev, and Andrew Semenov (2020). “QRE in the Ising game”. In:
3158 *CEUR Workshop*.
- 3159 — (2024). “Ising game on graphs”. In: *Chaos, Solitons & Fractals*.
- 3160 Leshno, Moshe, Vladimir Lin, Allan Pinkus, and Shimon Schocken (1993). “Multilayer feedforward
3161 networks with a nonpolynomial activation function can approximate any function”. In: *Neural
3162 networks*.
- 3163 Leslie, David S. and Edmund J. Collins (2003). “Convergent multiple-timescales reinforcement
3164 learning algorithms in normal form games”. In: *The Annals of Applied Probability*.
- 3165 — (2006). “Generalised weakened fictitious play”. In: *Games and Economic Behavior (GEB)*.
- 3166 Letcher, Alistair (2018). “Stability and exploitation in differentiable games”. MA thesis. University
3167 of Oxford.
- 3168 Letcher, Alistair, David Balduzzi, Sébastien Racanière, James Martens, Jakob Foerster, Karl Tuyls,
3169 and Thore Graepel (2019). “Differentiable game mechanics”. In: *Journal of Machine Learning
3170 Research (JMLR)*.
- 3171 Letcher, Alistair, Jakob Foerster, David Balduzzi, Tim Rocktäschel, and Shimon Whiteson (2019).
3172 “Stable opponent shaping in differentiable games”. In: *International Conference on Learning
3173 Representations (ICLR)*.
- 3174 Levine, David K. (1998). “Modeling altruism and spitefulness in experiments”. In: *Review of Economic
3175 Dynamics*.
- 3176 Li, Zhizhong and Derek Hoiem (2017). “Learning without forgetting”. In: *IEEE transactions on
3177 pattern analysis and machine intelligence*.
- 3178 Li, Zun and Michael P. Wellman (2021). “Evolution strategies for approximate solution of Bayesian
3179 games”. In: *AAAI Conference on Artificial Intelligence*.
- 3180 Lin, Henry W., Max Tegmark, and David Rolnick (2017). “Why does deep and cheap learning work
3181 so well?” In: *Journal of Statistical Physics*.
- 3182 Lisý, Viliam and Michael Bowling (2017). “Equilibrium approximation quality of current no-limit
3183 poker bots”. In: *AAAI Computer Poker Workshop*.

- 3184 Lockhart, Edward, Marc Lanctot, Julien Pérolat, Jean-Baptiste Lespiau, Dustin Morrill, Finbarr
3185 Timbers, and Karl Tuyls (2019). “Computing approximate equilibria in sequential adversarial
3186 games by exploitability descent”. In: *International Joint Conference on Artificial Intelligence*
3187 (*IJCAI*).
- 3188 Loizou, Nicolas, Hugo Berard, Gauthier Gidel, Ioannis Mitliagkas, and Simon Lacoste-Julien (2021).
3189 “Stochastic gradient descent-ascent and consensus optimization for smooth games: convergence
3190 analysis under expected co-coercivity”. In: *Conference on Neural Information Processing Systems*
3191 (*NeurIPS*).
- 3192 Lorraine, Jonathan and David Duvenaud (2018). “Stochastic hyperparameter optimization through
3193 hypernetworks”. In: *arXiv:1802.09419*.
- 3194 Lorraine, Jonathan, Paul Vicol, and David Duvenaud (2020). “Optimizing millions of hyperparameters
3195 by implicit differentiation”. In: *International conference on artificial intelligence and statistics*.
3196 Proceedings of Machine Learning Research (PMLR).
- 3197 Lotker, Zvi, Boaz Patt-Shamir, and Mark R. Tuttle (2008). “A game of timing and visibility”. In:
3198 *Games and Economic Behavior (GEB)*.
- 3199 Lu, Zhou, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang (2017). “The expressive power
3200 of neural networks: a view from the width”. In: *Conference on Neural Information Processing*
3201 *Systems (NeurIPS)*.
- 3202 Luce, R Duncan (1959). *Individual choice behavior*. Wiley New York.
- 3203 — (1977). “The choice axiom after twenty years”. In: *Journal of mathematical psychology*.
- 3204 Ma, Jeffrey, Alistair Letcher, Florian Schäfer, Yuanyuan Shi, and Anima Anandkumar (2021).
3205 “Polymatrix competitive gradient descent”. In: *arXiv:2111.08565*.
- 3206 Macal, C. M. and M. J. North (2010). “Tutorial on agent-based modelling and simulation”. In:
3207 *Journal of Simulation*.
- 3208 MacKay, Matthew, Paul Vicol, Jon Lorraine, David Duvenaud, and Roger Grosse (2019). “Self-tuning
3209 networks: bilevel optimization of hyperparameters using structured best-response functions”. In:
3210 *arXiv:1903.03088*.
- 3211 Marchesi, Alberto, Francesco Trovò, and Nicola Gatti (2020). “Learning probably approximately cor-
3212 rect maximin strategies in simulation-based games with infinite strategy spaces”. In: *International*
3213 *Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- 3214 Marris, Luke, Paul Muller, Marc Lanctot, Karl Tuyls, and Thore Graepel (2021). “Multi-agent
3215 training beyond zero-sum with correlated equilibrium meta-solvers”. In: *International Conference*
3216 *on Machine Learning (ICML)*.
- 3217 Martin, Carlos, Craig Boutilier, Ofer Meshi, and Tuomas Sandholm (2024). “Model-free preference
3218 elicitation”. In: *International Joint Conference on Artificial Intelligence (IJCAI)*.
- 3219 Martin, Carlos and Tuomas Sandholm (2021a). “Bayesian multiagent inverse reinforcement learn-
3220 ing for policy recommendation”. In: *AAAI Conference on Artificial Intelligence Workshop on*
3221 *Reinforcement Learning in Games*.
- 3222 — (2021b). “Efficient exploration of zero-sum stochastic games”. In: *AAAI Conference on Artificial*
3223 *Intelligence Workshop on Reinforcement Learning in Games*.
- 3224 — (2023). “Finding mixed-strategy equilibria of continuous-action games without gradients using
3225 randomized policy networks”. In: *International Joint Conference on Artificial Intelligence (IJCAI)*.
- 3226 — (2025a). “AlphaZeroES: direct score maximization outperforms planning loss minimization”. In:
3227 *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- 3228 — (2025b). “ApproxED: approximate exploitability descent via learned best responses”. In: *Interna-*
3229 *tional Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.

3230 Martin, Carlos and Tuomas Sandholm (2025c). “Incremental multiple oracle”. In: *International*
3231 *Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
3232 — (2025d). “Joint-perturbation simultaneous pseudo-gradient”. In: *International Joint Conference*
3233 *on Artificial Intelligence (IJCAI)*.
3234 — (2025e). “Solving infinite-player games with player-to-strategy networks”. In: *arXiv:2501.09330*.
3235 Mazumdar, Eric, Lillian J. Ratliff, and S. Shankar Sastry (2020). “On gradient-based learning in
3236 continuous games”. In: *SIAM Journal on Mathematics of Data Science (SIMODS)*.
3237 Mazumdar, Eric, S. Shankar Sastry, and Michael I. Jordan (2025). “On finding local Nash equilibria
3238 (and only local Nash equilibria) in zero-sum games”. In: *ACM/IMS Journal of Data Science*.
3239 McAleer, Stephen, JB Lanier, Kevin A. Wang, Pierre Baldi, Tuomas Sandholm, and Roy Fox (2024).
3240 “Toward optimal policy population growth in two-player zero-sum games”. In: *International*
3241 *Conference on Learning Representations (ICLR)*.
3242 McAleer, Stephen, John B. Lanier, Roy Fox, and Pierre Baldi (2020). “Pipeline PSRO: a scalable
3243 approach for finding approximate Nash equilibria in large games”. In: *Conference on Neural*
3244 *Information Processing Systems (NeurIPS)*.
3245 McAleer, Stephen, John B. Lanier, Kevin A. Wang, Pierre Baldi, and Roy Fox (2021). “XDO:
3246 a double oracle algorithm for extensive-form games”. In: *Conference on Neural Information*
3247 *Processing Systems (NeurIPS)*.
3248 McAleer, Stephen, Kevin Wang, John Lanier, Marc Lanctot, Pierre Baldi, Tuomas Sandholm, and
3249 Roy Fox (2022). “Anytime PSRO for two-player zero-sum games”. In: *arXiv:2201.07700*.
3250 McCloskey, Michael and Neal J. Cohen (1989). “Catastrophic interference in connectionist networks:
3251 the sequential learning problem”. In: *Psychology of learning and motivation*.
3252 McCulloch, Warren S. and Walter Pitts (1943). “A logical calculus of the ideas immanent in nervous
3253 activity”. In: *The bulletin of mathematical biophysics*.
3254 McKee, Kevin R., Ian Gemp, Brian McWilliams, Edgar A. Duñez-Guzmán, Edward Hughes, and Joel
3255 Z. Leibo (2020). “Social diversity and social preferences in mixed-motive reinforcement learning”.
3256 In: *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
3257 McKelvey, Richard D. and Thomas R. Palfrey (1995). “Quantal response equilibria for normal form
3258 games”. In: *Games and economic behavior*.
3259 — (1998). “Quantal response equilibria for extensive form games”. In: *Experimental economics*.
3260 McKinney, Wes (2010). “Data structures for statistical computing in Python”. In: *Proceedings of the*
3261 *9th Python in Science Conference*. Ed. by Stéfan van der Walt and Jarrod Millman.
3262 McMahan, H. Brendan, Geoffrey J. Gordon, and Avrim Blum (2003). “Planning in the presence of
3263 cost functions controlled by an adversary”. In: *International Conference on Machine Learning*
3264 *(ICML)*.
3265 McNamara, John M. and Alasdair I. Houston (1996). “State-dependent life histories”. In: *Nature*.
3266 Mertikopoulos, Panayotis and Mathias Staudigl (2018). “On the convergence of gradient-like flows
3267 with noisy gradient input”. In: *SIAM Journal on Optimization*.
3268 Mertikopoulos, Panayotis and Zhengyuan Zhou (2019). “Learning in games with continuous action
3269 sets and unknown payoff functions”. In: *Mathematical Programming*.
3270 Merz, Antony W. (1974). “The homicidal chauffeur”. In: *AIAA Journal*.
3271 Mescheder, Lars, Sebastian Nowozin, and Andreas Geiger (2017). “The numerics of GANs”. In:
3272 *Conference on Neural Information Processing Systems (NeurIPS)*.
3273 Metropolis, Nicholas, Arianna Rosenbluth, Marshall Rosenbluth, Augusta Teller, and Edward Teller
3274 (1953). “Equation of state calculations by fast computing machines”. In: *The Journal of Chemical*
3275 *Physics*.

- 3276 Metz, Luke, Ben Poole, David Pfau, and Jascha Sohl-Dickstein (2017). “Unrolled generative adver-
3277 sarial networks”. In: *International Conference on Learning Representations (ICLR)*.
- 3278 Milchtaich, Igal (1996). *Generic uniqueness of equilibria in nonatomic congestion games*. Hebrew
3279 University of Jerusalem.
- 3280 — (2000). “Generic uniqueness of equilibrium in large crowding games”. In: *Mathematics of Opera-*
3281 *tions Research*.
- 3282 — (2005). “Topological conditions for uniqueness of equilibrium in networks”. In: *Mathematics of*
3283 *Operations Research*.
- 3284 Milgrom, Paul and Ilya Segal (2014). “Deferred-acceptance auctions and radio spectrum reallocation”.
3285 In: *ACM Conference on Economics and Computation (EC)*.
- 3286 — (2020). “Clock auctions and radio spectrum reallocation”. In: *Journal of Political Economy*.
- 3287 Milgrom, Paul and Robert Weber (1985). “Distributional strategies for games with incomplete
3288 information”. In: *Mathematics of Operations Research*.
- 3289 Mitchell, Tom (1997). *Machine Learning*. McGraw-Hill.
- 3290 Mobiny, Aryan, Pengyu Yuan, Supratik K. Moulik, Naveen Garg, Carol C. Wu, and Hien Van
3291 Nguyen (2021). “DropConnect is effective in modeling uncertainty of Bayesian deep networks”.
3292 In: *Scientific reports*.
- 3293 Monderer, Dov and Lloyd S. Shapley (1996). “Potential games”. In: *Games and Economic Behavior*
3294 *(GEB)*.
- 3295 Moravčík, Matej, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis,
3296 Kevin Waugh, Michael Johanson, and Michael Bowling (2017). “DeepStack: expert-level artificial
3297 intelligence in heads-up no-limit poker”. In: *Science*.
- 3298 Morgan, John, Kenneth Steiglitz, and George Reis (2003). “The spite motive and equilibrium
3299 behavior in auctions”. In: *Contributions to Economic Analysis & Policy*.
- 3300 Moulin, H. and J.-P. Vial (1978). “Strategically zero-sum games: the class of games whose completely
3301 mixed equilibria cannot be improved upon”. In: *International Journal of Game Theory*.
- 3302 Muller, Paul, Shayegan Omidshafiei, Mark Rowland, Karl Tuyls, Julien Perolat, Siqi Liu, Daniel
3303 Hennes, Luke Marris, Marc Lanctot, Edward Hughes, Zhe Wang, Guy Lever, Nicolas Heess, Thore
3304 Graepel, and Remi Munos (2020). “A generalized training approach for multiagent learning”. In:
3305 *International Conference on Learning Representations (ICLR)*.
- 3306 Muller, Paul, Mark Rowland, Romuald Elie, Georgios Piliouras, Julien Perolat, Mathieu Lauriere,
3307 Raphael Marinier, Olivier Pietquin, and Karl Tuyls (2022). “Learning equilibria in mean-field
3308 games: introducing mean-field PSRO”. In: *International Conference on Autonomous Agents and*
3309 *Multi-Agent Systems (AAMAS)*.
- 3310 Myerson, Roger B. (1991). *Game theory: analysis of Conflict*. Harvard University Press.
- 3311 Nash, John F. (1950). “Equilibrium points in n-person games”. In: *Proceedings of the National*
3312 *Academy of Sciences (PNAS)*.
- 3313 — (1951). “Non-cooperative games”. In: *Annals of Mathematics*.
- 3314 Neal, Radford M. (1996). “Monte Carlo implementation”. In: *Bayesian learning for neural networks*.
3315 Springer.
- 3316 Nesterov, Yurii and Vladimir Spokoiny (2017). “Random gradient-free minimization of convex
3317 functions”. In: *Foundations of Computational Mathematics*.
- 3318 Niederreiter, Harald (1992). *Random number generation and quasi-Monte Carlo methods*. SIAM.
- 3319 Nikaido, Hukukane and Kazuo Isoda (1955). “Note on non-cooperative convex games”. In: *Pacific*
3320 *Journal of Mathematics*.

- 3321 Novshek, William (1985). “Perfectly competitive markets as the limits of Cournot markets”. In:
3322 *Journal of Economic Theory (JET)*.
- 3323 Omidshafiei, Shayegan, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland,
3324 Jean-Baptiste Lespiau, Wojciech Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos
3325 (2019). “ α -rank: multi-agent evaluation by evolution”. In: *Scientific Reports*.
- 3326 Osborne, Martin J. (2004). *An introduction to game theory*. Oxford University Press.
- 3327 Padala, Manisha, Debojit Das, and Sujit Gujar (2021). “Effect of input noise dimension in GANs”.
3328 In: *Neural Information Processing*.
- 3329 Parise, Francesca and Asuman Ozdaglar (2019). “Graphon games”. In: *Economics and Computation*.
3330 — (2023). “Graphon games: a statistical framework for network games and interventions”. In:
3331 *Econometrica*.
- 3332 Patsko, Valery S. and Varvara L. Turova (2001). “Level sets of the value function in differential
3333 games with the homicidal chauffeur dynamics”. In: *International Game Theory Review*.
- 3334 Perkins, Steven and David S. Leslie (2014). “Stochastic fictitious play with continuous action sets”.
3335 In: *Journal of Economic Theory (JET)*.
- 3336 Perolat, Julien, Bart De Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer,
3337 Paul Muller, Jerome T. Connor, Neil Burch, Thomas Anthony, Stephen McAleer, Romuald Elie,
3338 Sarah H. Cen, Zhe Wang, Audrunas Gruslys, Aleksandra Malysheva, Mina Khan, Sherjil Ozair,
3339 Finbarr Timbers, Toby Pohlen, Tom Eccles, Mark Rowland, Marc Lanctot, Jean-Baptiste Lespiau,
3340 Bilal Piot, Shayegan Omidshafiei, Edward Lockhart, Laurent Sifre, Nathalie Beauguerlange,
3341 Remi Munos, David Silver, Satinder Singh, Demis Hassabis, and Karl Tuyls (2022). “Mastering
3342 the game of Stratego with model-free multiagent reinforcement learning”. In: *Science*.
- 3343 Pérolat, Julien, Sarah Perrin, Romuald Elie, Mathieu Laurière, Georgios Piliouras, Matthieu Geist,
3344 Karl Tuyls, and Olivier Pietquin (2022). “Scaling mean field games by online mirror descent”. In:
3345 *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- 3346 Perrin, Sarah, Julien Perolat, Mathieu Laurière, Matthieu Geist, Romuald Elie, and Olivier Pietquin
3347 (2020). “Fictitious play for mean field games: continuous time analysis and applications”. In:
3348 *Conference on Neural Information Processing Systems (NeurIPS)*.
- 3349 Petersen, Peter (2006). *Riemannian geometry*. Springer.
- 3350 Pettis, Billy James (1938). “On integration in vector spaces”. In: *Transactions of the American
3351 Mathematical Society*.
- 3352 Pillards, Tim and Ronald Cools (2005). “Transforming low-discrepancy sequences from a cube to a
3353 simplex”. In: *Journal of computational and applied mathematics*.
- 3354 Pinkus, Allan (1999). “Approximation theory of the MLP model in neural networks”. In: *Acta
3355 numerica*.
- 3356 Pitis, Silviu (2019). “Rethinking the discount factor in reinforcement learning: a decision theoretic
3357 approach”. In: *AAAI Conference on Artificial Intelligence*.
- 3358 Plackett, Robin L. (1975). “The analysis of permutations”. In: *Journal of the Royal Statistical Society
3359 Series C: Applied Statistics*.
- 3360 Poisson, Siméon-Denis (1837). *Recherches sur la probabilité des jugements en matière criminelle et
3361 en matière civile: précédées des règles générales du calcul des probabilités*. Bachelier.
- 3362 Popov, Leonid (1980). “A modification of the Arrow–Hurwicz method for search of saddle points”.
3363 In: *Mathematical Notes*.
- 3364 Qin, Rong-Jun, Fan-Ming Luo, Hong Qian, and Yang Yu (2022). “Unified policy optimization for
3365 continuous-action reinforcement learning in non-stationary tasks and games”. In: *arXiv:2208.09452*.

- 3366 Qu, Biao and Jing Zhao (2013). “Methods for solving generalized Nash equilibrium”. In: *Journal of*
3367 *Applied Mathematics*.
- 3368 Raghunathan, Arvind, Anoop Cherian, and Devesh Jha (2019). “Game theoretic optimization via
3369 gradient-based Nikaido–Isoda function”. In: *International Conference on Machine Learning*
3370 *(ICML)*.
- 3371 Rahaman, Nasim, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua
3372 Bengio, and Aaron Courville (2019). “On the spectral bias of neural networks”. In: *International*
3373 *Conference on Machine Learning (ICML)*.
- 3374 Rath, Kali (1992). “A direct proof of the existence of pure strategy equilibria in games with a
3375 continuum of players”. In: *Economic Theory*.
- 3376 Ray, Debraj and Rajiv Vohra (2020). “Games of love and hate”. In: *Journal of Political Economy*.
- 3377 Rebuffi, Sylvestre-Alvise, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert (2017).
3378 “iCaRL: incremental classifier and representation learning”. In: *IEEE Conference on Computer*
3379 *Vision and Pattern Recognition (CVPR)*.
- 3380 Rechenberg, Ingo (1973). *Evolutionsstrategie: optimierung technischer systeme nach prinzipien der*
3381 *biologischen evolution*. Frommann-Holzboog.
- 3382 — (1978). “Evolutionsstrategien”. In: *Simulationsmethoden in der medizin und biologje*. Springer.
- 3383 Reluga, Timothy C. (2010). “Game theory of social distancing in response to an epidemic”. In: *PLoS*
3384 *Computational Biology*.
- 3385 Rhoads, Glenn C. and Laurent Bartholdi (2012). “Computer solution to the game of pure strategy”.
3386 In: *Games*.
- 3387 Riemann, Bernhard (1868). “Über die Darstellbarkeit einer Function durch eine trigonometrische
3388 Reihe”. In: *Abhandlungen der Königlichen Gesellschaft der Wissenschaften in Göttingen*.
- 3389 Roberts, Martin (2018). *The unreasonable effectiveness of quasirandom sequences*.
- 3390 Rosen, J. Ben (1965). “Existence and uniqueness of equilibrium points for concave n-person games”.
3391 In: *Econometrica*.
- 3392 Rosenblatt, Frank (1958). “The perceptron: a probabilistic model for information storage and
3393 organization in the brain”. In: *Psychological review*.
- 3394 Ross, Sheldon M. (1971). “Goofspiel: the game of pure strategy”. In: *Journal of Applied Probability*.
- 3395 Rossum, Guido Van and Fred Drake Jr (1995). *Python reference manual*. Centrum voor Wiskunde
3396 en Informatica Amsterdam.
- 3397 Rusu, Andrei A., Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick,
3398 Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell (2016). “Progressive neural networks”.
3399 In: *arXiv:1606.04671*.
- 3400 Saijo, Tatsuyoshi and Hideki Nakamura (1995). “The “spite” dilemma in voluntary contribution
3401 mechanism experiments”. In: *Journal of Conflict Resolution*.
- 3402 Saks, Stanisław (1937). *Theory of the Integral*. Monografie matematyczne.
- 3403 Salimans, Tim, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever (2017). “Evolution strategies
3404 as a scalable alternative to reinforcement learning”. In: *arXiv:1703.03864*.
- 3405 Sandholm, Tuomas (2013). “Very-large-scale generalized combinatorial multi-attribute auctions:
3406 lessons from conducting \$60 billion of sourcing”. In: *Handbook of Market Design*. Oxford University
3407 Press.
- 3408 Sandholm, William H. (2001). “Potential games with continuous player sets”. In: *Journal of Economic*
3409 *Theory*.
- 3410 Santos, Manuel S. and Michael Woodford (1997). “Rational asset pricing bubbles”. In: *Econometrica:*
3411 *Journal of the Econometric Society*.

- 3412 Sard, Arthur (1942). “The measure of the critical values of differentiable maps”. In: *Bulletin of the*
3413 *American Mathematical Society*.
- 3414 Schäfer, Florian and Anima Anandkumar (2019). “Competitive gradient descent”. In: *Conference on*
3415 *Neural Information Processing Systems (NeurIPS)*.
- 3416 Schmeidler, David (1973). “Equilibrium points of nonatomic games”. In: *Journal of Statistical Physics*.
- 3417 Schmidhuber, Jürgen (1992). “Learning to control fast-weight memories: an alternative to dynamic
3418 recurrent networks”. In: *Neural Computation*.
- 3419 Schwarz, Hermann Amandus (1890). “Ueber ein die Flächen kleinsten Flächeninhalts betreffendes
3420 Problem der Variationsrechnung: Festschrift zum siebzigsten Geburtstage des Herrn Karl Weier-
3421 strass”. In: *Gesammelte Mathematische Abhandlungen: Erster Band*. Springer.
- 3422 Schwefel, Hans-Paul (1977). *Numerische optimierung von computer-modellen mittels der evolution-
3423 sstrategie*. Birkhäuser Basel.
- 3424 Selten, Reinhard (1970). *Preispolitik der Mehrproduktenunternehmung in der statischen Theorie*.
3425 Springer Berlin, Heidelberg.
- 3426 Shamir, Ohad (2017). “An optimal algorithm for bandit and zero-order convex optimization with
3427 two-point feedback”. In: *Journal of Machine Learning Research (JMLR)*.
- 3428 Shapley, Lloyd S. (1964). “Some topics in two-person games”. In: *Advances in Game Theory*. Princeton
3429 University Press.
- 3430 Shapley, Lloyd S. and Martin Shubik (1971). “The assignment game I: the core”. In: *International*
3431 *Journal of Game Theory*.
- 3432 Sharma, Ankit and Tuomas Sandholm (2010). “Asymmetric spite in auctions”. In: *AAAI Conference*
3433 *on Artificial Intelligence*.
- 3434 Shin, Hanul, Jung Kwon Lee, Jaehong Kim, and Jihoon Kim (2017). “Continual learning with deep
3435 generative replay”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 3436 Silver, David, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,
3437 Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen
3438 Simonyan, and Demis Hassabis (2018). “A general reinforcement learning algorithm that masters
3439 chess, shogi, and Go through self-play”. In: *Science*.
- 3440 Simon, Herbert A. (1955). “A behavioral model of rational choice”. In: *The Quarterly Journal of*
3441 *Economics*.
- 3442 Sinha, Arunesh, Fei Fang, Bo An, Christopher Kiekintveld, and Milind Tambe (2018). “Stackelberg
3443 security games: looking beyond a decade of success”. In: *International Joint Conference on*
3444 *Artificial Intelligence (IJCAI)*.
- 3445 Sion, Maurice (1958). “On general minimax theorems”. In: *Pacific Journal of Mathematics*.
- 3446 Smith, John Maynard (1974). “The theory of games and evolution in animal conflict”. In: *Journal of*
3447 *Theoretical Biology*.
- 3448 Spall, James C. (1992). “Multivariate stochastic approximation using a simultaneous perturbation
3449 gradient approximation”. In: *IEEE Transactions on Automatic Control*.
- 3450 — (1997). “A one-measurement form of simultaneous perturbation stochastic approximation”. In:
3451 *Automatica*.
- 3452 Srivastava, Nitish, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov
3453 (2014). “Dropout: a simple way to prevent neural networks from overfitting”. In: *Journal of*
3454 *Machine Learning Research (JMLR)*.
- 3455 Steinberger, Eric, Adam Lerer, and Noam Brown (2020). “DREAM: deep regret minimization with
3456 advantage baselines and model-free learning”. In: *arXiv:2006.10410*.

- 3457 Stieltjes, Thomas Joannes (1894). “Recherches sur les fractions continues”. In: *Annales de la Faculté*
3458 *des sciences de Toulouse: Mathématiques*.
- 3459 Sutton, Richard S. and Andrew G. Barto (1998). *Reinforcement learning: an introduction*. MIT
3460 press Cambridge.
- 3461 Szafron, Duane, Richard G. Gibson, and Nathan R. Sturtevant (2013). “A parameterized family of
3462 equilibrium profiles for three-player Kuhn poker”. In: *International Conference on Autonomous*
3463 *Agents and Multi-Agent Systems (AAMAS)*.
- 3464 Szentes, Balázs and Robert W. Rosenthal (2003a). “Beyond chopsticks: symmetric equilibria in
3465 majority auction games”. In: *Games and Economic Behavior*.
- 3466 — (2003b). “Three-object two-bidder simultaneous auctions: chopsticks and tetrahedra”. In: *Games*
3467 *and Economic Behavior (GEB)*.
- 3468 Talvila, Erik (2001). “Necessary and sufficient conditions for differentiating under the integral sign”.
3469 In: *The American Mathematical Monthly*.
- 3470 Tancik, Matthew, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh
3471 Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng (2020). “Fourier features let networks
3472 learn high frequency functions in low dimensional domains”. In: *Conference on Neural Information*
3473 *Processing Systems (NeurIPS)*.
- 3474 Tang, Pingzhong and Tuomas Sandholm (2012). “Optimal auctions for spiteful bidders”. In: *AAAI*
3475 *Conference on Artificial Intelligence*.
- 3476 Telgarsky, Matus (2016). “Benefits of depth in neural networks”. In: *Conference on Learning Theory*
3477 *(COLT)*.
- 3478 The HDF Group (2024). *Hierarchical data format, version 5*. DOI: 10.5281/zenodo.17808614.
- 3479 The Pandas development team (2020). *Pandas*. DOI: 10.5281/zenodo.3509134.
- 3480 Timbers, Finbarr, Nolan Bard, Edward Lockhart, Marc Lanctot, Martin Schmid, Neil Burch, Julian
3481 Schrittwieser, Thomas Hubert, and Michael Bowling (2022). “Approximate exploitability: learning
3482 a best response”. In: *International Joint Conference on Artificial Intelligence (IJCAI)*.
- 3483 Tsaknakis, Ioannis and Mingyi Hong (2021). “Finding first-order Nash equilibria of zero-sum
3484 games with the regularized Nikaido–Isoda function”. In: *International Conference on Artificial*
3485 *Intelligence and Statistics (AISTATS)*.
- 3486 Tucker, Albert W. (1984). *The Princeton mathematics community in the 1930s: transcript number*
3487 *11*.
- 3488 Uryasev, Stanislav and Reuven Y. Rubinstein (1994). “On relaxation algorithms in computation of
3489 noncooperative equilibria”. In: *IEEE Transactions on Automatic Control (TACON)*.
- 3490 Vickrey, William S. (1969). “Congestion theory and transport investment”. In: *American Economic*
3491 *Review (AER)*.
- 3492 Vinyals, Oriol, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung
3493 Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan,
3494 Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max
3495 Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David
3496 Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff,
3497 Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith,
3498 Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David
3499 Silver (2019). “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. In:
3500 *Nature*.
- 3501 Virtanen, Pauli, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournau-
3502 peau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der

- 3503 Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson,
3504 Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake
3505 VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero,
3506 Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt,
3507 and SciPy 1.0 Contributors (2020). “SciPy 1.0: fundamental Algorithms for Scientific Computing
3508 in Python”. In: *Nature Methods*.
- 3509 Vitter, Jeffrey S. (1985). “Random sampling with a reservoir”. In: *ACM Transactions on Mathematical
3510 Software (TOMS)*.
- 3511 Vlatakis-Gkaragkounis, Emmanouil-Vasileios, Lampros Flokas, Thanasis Lianeas, Panayotis Mer-
3512 tikopoulos, and Georgios Piliouras (2020). “No-regret learning and mixed nash equilibria: They
3513 do not mix”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 3514 von Neumann, John (1928). “Zur theorie der gesellschaftsspiele”. In: *Mathematische Annalen*.
- 3515 von Neumann, John and Oskar Morgenstern (1947). *Theory of games and economic behavior*.
3516 Princeton University Press.
- 3517 Voorhees, Alan M. (1955). “A general theory of traffic movement”. In: *Institute of Traffic Engineers*.
- 3518 Walton, Michael and Viliam Lisy (2021). “Multi-agent reinforcement learning in OpenSpiel: a
3519 reproduction report”. In: *arXiv:2103.00187*.
- 3520 Wan, Li, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus (2013). “Regularization of
3521 neural networks using DropConnect”. In: *International Conference on Machine Learning (ICML)*.
3522 Proceedings of Machine Learning Research (PMLR).
- 3523 Wang, Yongzhao and Michael Wellman (2023). “Empirical game-theoretic analysis for mean field
3524 games”. In: *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- 3525 Wang, Yuanhao, Guodong Zhang, and Jimmy Ba (2020). “On solving minimax optimization locally:
3526 a follow-the-ridge approach”. In: *International Conference on Learning Representations (ICLR)*.
- 3527 Wardrop, John Glen (1952). “Some theoretical aspects of road traffic research”. In: *Proceedings of
3528 the Institution of Civil Engineers*.
- 3529 Washburn, Alan (2013). “OR forum - Blotto politics”. In: *Operations Research*.
- 3530 Wellman, Michael P., Karl Tuyls, and Amy Greenwald (2025). “Empirical game theoretic analysis: a
3531 survey”. In: *Journal of Artificial Intelligence Research (JAIR)*.
- 3532 White, Martha (2017). “Unifying task specification in reinforcement learning”. In: *International
3533 Conference on Machine Learning (ICML)*.
- 3534 Wiener, Norbert (1923). “Differential-space”. In: *Journal of Mathematics and Physics*.
- 3535 — (1938). “The homogeneous chaos”. In: *American Journal of Mathematics*.
- 3536 Wierstra, Daan, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber
3537 (2014). “Natural evolution strategies”. In: *Journal of Machine Learning Research (JMLR)*.
- 3538 Willi, Timon, Alistair Letcher, Johannes Treutlein, and Jakob Foerster (2022). “COLA: consistent
3539 learning with opponent-learning awareness”. In: *International Conference on Machine Learning
3540 (ICML)*.
- 3541 Wilson, Alan G. (1967). “A statistical theory of spatial distribution models”. In: *Transportation
3542 Research*.
- 3543 Wiszniewska-Matyskiel, Agnieszka (2014). “Open and closed loop Nash equilibria in games with a
3544 continuum of players”. In: *Journal of Optimization Theory and Applications*.
- 3545 Worth, Carl and Keith Packard (2003). “Xr: cross-device rendering for vector graphics”. In: *Proceedings
3546 of the Ottawa Linux Symposium (OLS)*.
- 3547 Wu, Zida, Mathieu Laurière, Samuel Jia Cong Chua, Matthieu Geist, Olivier Pietquin, and Ankur
3548 Mehta (2024). “Population-aware Online Mirror Descent for Mean-Field Games by Deep Re-

- 3549 reinforcement Learning”. In: *International Conference on Autonomous Agents and Multi-Agent*
3550 *Systems (AAMAS)*.
- 3551 Xie, Qiaomin, Zhuoran Yang, Zhaoran Wang, and Andreea Minca (2021). “Learning while playing
3552 in mean-field games: Convergence and optimality”. In: *International Conference on Machine*
3553 *Learning (ICML)*.
- 3554 Xin, Chen, G. Yang, and J. P. Huang (2017). “Ising game: nonequilibrium steady states of resource-
3555 allocation systems”. In: *Physica A*.
- 3556 Yang, Zhe and Qingping Song (2022). “A unified approach to the Nash equilibrium existence in
3557 large games from finitely many players to infinitely many players”. In: *Journal of Fixed Point*
3558 *Theory and Applications*.
- 3559 Yarotsky, Dmitry (2017). “Error bounds for approximations with deep ReLU networks”. In: *Neural*
3560 *Networks*.
- 3561 Yellott Jr, John I. (1977). “The relationship between Luce’s choice axiom, Thurstone’s theory of
3562 comparative judgment, and the double exponential distribution”. In: *Journal of Mathematical*
3563 *Psychology*.
- 3564 Yi, Sun, Daan Wierstra, Tom Schaul, and Jürgen Schmidhuber (2009). “Stochastic search using the
3565 natural gradient”. In: *International Conference on Machine Learning (ICML)*.
- 3566 Yosida, Kôsaku and Edwin Hewitt (1952). “Finitely additive measures”. In: *Transactions of the*
3567 *American Mathematical Society*.
- 3568 Zhang, Chenyu, Xu Chen, and Xuan Di (2025). “Stochastic semi-gradient descent for learning mean
3569 field games with population-aware function approximation”. In: *International Conference on*
3570 *Learning Representations (ICLR)*.
- 3571 Zhang, Chongjie and Victor Lesser (2010). “Multi-agent learning with policy prediction”. In: *AAAI*
3572 *Conference on Artificial Intelligence*.
- 3573 Zhuang, Juntang, Tommy Tang, Yifan Ding, Sekhar C. Tatikonda, Nicha Dvornek, Xenophon
3574 Papademetris, and James Duncan (2020). “AdaBelief optimizer: adapting stepsizes by the belief
3575 in observed gradients”. In: *Conference on Neural Information Processing Systems (NeurIPS)*.
- 3576 Zinkevich, Martin, Michael Bowling, Michael Johanson, and Carmelo Piccione (2007). “Regret
3577 minimization in games with incomplete information”. In: *Conference on Neural Information*
3578 *Processing Systems (NeurIPS)*.

3579 Appendix A

3580 Other completed work

3581 In addition to the work presented in this thesis, I have also completed other projects during my PhD
3582 that, while valuable, fall outside the primary focus of this thesis project. These efforts demonstrate
3583 my broader research capabilities and collaborative experience.

- 3584 • In Martin and Sandholm (2021b), we studied **efficient exploration of two-player zero-sum**
3585 **stochastic games** when the learner can control both players, only has oracle access via gameplay,
3586 and has a limited-duration exploration phase during which it can control both players. The
3587 objective is to output a strategy with low exploitability after exploration. We proposed using a
3588 belief distribution over possible environments to induce a distribution over state–action value
3589 functions and evaluate exploration strategies (including generalizations of Thompson sampling
3590 and Bayes-UCB), finding that Thompson-style and Bayes-UCB–style methods consistently
3591 outperform alternatives.
- 3592 • In Martin and Sandholm (2021a), we studied **Bayesian multiagent inverse reinforcement**
3593 **learning for policy recommendation**. In this setting, an observer (the recommender) watches
3594 players’ actions and state transitions in a known zero-sum game but not their rewards, and uses
3595 a Bayesian posterior (with behavioral models that may be Nash, quantal response, etc.) to infer
3596 players’ reward functions. We proposed Bayesian inference procedures to recover posteriors over
3597 rewards/behaviors and show that the inferred preferences can be used to distinguish solution
3598 concepts and drive policy recommendations.
- 3599 • In Martin, Boutilier, et al. (2024), we tackled **model-free preference elicitation** in recom-
3600 mender systems. Preference elicitation allows the system to ask a user questions to learn about
3601 their preferences and improve recommendation quality. We developed model-free variants of
3602 **expected value of information (EVOI)** that learn user-response and utility models from
3603 existing data (using function approximation) so preference-elicitation queries can be scored
3604 without explicit probabilistic preference models, and augment these learned models with online
3605 planning via Monte-Carlo tree search. Empirically, this approach yielded significant improvements
3606 in recommendation quality over standard baselines across several preference-elicitation tasks.
- 3607 • Planning at execution time can dramatically improve performance in some settings. AlphaZero
3608 (Silver et al., 2018) is a state-of-the-art technique based on Monte Carlo Tree Search (MCTS)
3609 (Coulom, 2007). In Martin and Sandholm (2025a), we introduced **AlphaZeroES**. It modifies
3610 AlphaZero by replacing its loss with direct maximization of the episode score, while keeping

3611 the MCTS algorithm and neural architecture unchanged. Since the episode score is generally
3612 not differentiable, we used evolution strategies (as described in Section 2.8) for this. Across
3613 several single-agent environments, experiments showed that directly maximizing episode score
3614 consistently and often dramatically outperforms minimizing the standard planning loss.

- 3615 • In Martin and Sandholm (2025c), we presented **incremental multiple oracle (IMO)**, a
3616 framework for computing approximate mixed-strategy Nash equilibria of continuous-action
3617 games. It is a modification of the traditional double oracle algorithm, extended to multiple
3618 players and continuous action spaces. It does not require exact metagame solving on each
3619 iteration, which can be computationally expensive for large metagames. Also, it does not require
3620 global best-response computation on each iteration, which can be computationally expensive or
3621 even intractable for high-dimensional action spaces and general games.

3622 Appendix B

3623 Additional information

3624 Here, we include some additional information pertaining to the content presented in the body.

3625 B.1 Additional information about auctions

3626 Table B.1 describes the independent private values, common values, affiliated values, complete
3627 information, and asymmetric information auctions, in that order.

Ω	$\tau_i(\omega)$	$v_i(\omega)$
$[0, 1]^n$	ω_i	ω_i
$[0, 1]^{n+1}$	$\omega_i \omega_{n+1}$	ω_{n+1}
$[0, 1]^{n+1}$	$\omega_i + \omega_{n+1}$	$\omega_{n+1} + \frac{1}{n} \sum_{i=1}^n \omega_i$
$[0, 1]$	ω	ω
$[0, 1]$	$\omega \mathbb{I}[i = 1]$	ω

Table B.1: Auction descriptions. All use $\mu = \mathcal{U}(\Omega)$.

3628 For the common values auction, also known as the “mineral rights” model (Krishna, 2002, example
3629 6.1), the following procedure can be used to sample $\omega \mid o_i$:

$$z \sim \mathcal{U}([0, 1]) \tag{B.1}$$

$$\omega_{n+1} = o_i^z \tag{B.2}$$

$$\omega_i = o_i / \omega_{n+1} \tag{B.3}$$

$$\omega_j \sim \mathcal{U}([0, 1]) \quad j \neq i \tag{B.4}$$

3630 For the affiliated values auction (Krishna, 2002, example 6.2), the following procedure can be used
3631 to sample $\omega \mid o_i$:

$$\omega_{n+1} \sim \mathcal{U}(\max\{0, o_i - 1\}, \min\{1, o_i\}) \tag{B.5}$$

$$\omega_i = o_i - \omega_{n+1} \tag{B.6}$$

$$\omega_j \sim \mathcal{U}([0, 1]) \quad j \neq i \tag{B.7}$$

3632 The all-pay auction with independent private values has a pure symmetric equilibrium generated
 3633 by:

$$a_i = \frac{n-1}{n}(o_i)^n \tag{B.8}$$

3634 The k th-price winner-pay auction with independent private values has a pure symmetric equilibrium
 3635 generated by Kagel and Levin (1993, p. 878):

$$a_i = \frac{n-1}{n+1-k}o_i \tag{B.9}$$

3636 The 3-player 2nd-price winner-pay auction with common values has a pure symmetric equilibrium
 3637 generated by Bichler, Fichtl, Heidekrüger, et al. (2021):

$$a_i = \frac{o_i}{2+\frac{1}{2}o_i} \tag{B.10}$$

3638 The 2-player 1st- and 2nd-price winner-pay auction with affiliated values have pure symmetric
 3639 equilibria generated by, respectively (Bichler, Fichtl, Heidekrüger, et al., 2021):

$$a_i = \frac{2}{3}o_i \tag{B.11}$$

$$a_i = o_i \tag{B.12}$$

3640 B.2 Noise dimensionality in RPNs

3641 For randomized policy networks (Section 5.1), the dimensionality of the input noise is crucial to
 3642 performance. In this section, we review some results from the literature pertaining to this.

3643 As described in Section 1.1, universal approximation theorems guarantee that neural networks
 3644 can approximate arbitrary continuous functions on a compact domain.

3645 Furthermore, if $\mathcal{X}, \mathcal{Y} \subseteq \mathbb{R}^m$, X is an \mathcal{X} -valued random variable, $f : \mathcal{X} \rightarrow \mathcal{Y}$, and f_n is a sequence
 3646 of functions that converges pointwise to f , the sequence of random variables $Y_n = f_n(X)$ converges
 3647 in distribution to $Y = f(X)$ (Huang, Krueger, et al., 2018, Lemma 4).

3648 Gaussian input noise of *lower* dimension than the output space does not suffice to approximate
 3649 arbitrary distributions on the output space. In particular, **Sard’s theorem** (Sard, 1942; Petersen,
 3650 2006) says the following: Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be k times continuously differentiable, where $k \geq$
 3651 $\max\{n - m + 1, 1\}$. Let $X \subseteq \mathbb{R}^n$ be the set of critical points of f (that is, points where the Jacobian
 3652 matrix of f has rank less than m). Then the image $f[X]$ has Lebesgue measure zero in \mathbb{R}^m . As a
 3653 corollary, if $n < m$, then all points in \mathbb{R}^n are critical. Thus $\text{im } f$ has Lebesgue measure zero in \mathbb{R}^m .¹

3654 Padala, Das, and Gujar (2021) study the effect of input noise dimension in GANs. They show
 3655 that the right dimension of input noise for optimal results depends on the dataset and architecture
 3656 used.

3657 Feng, Zhao, and Zha (2021) study how noise injection in GANs helps them overcome the
 3658 “adversarial dimension trap”, which arises when the generated manifold has an intrinsic dimension
 3659 lower than that of the data manifold: that is, when the latent space is low-dimensional compared to
 3660 the high-dimensional space of real image details. Citing Sard’s theorem, they advise against mapping
 3661 low-dimensional feature spaces into feature manifolds with higher intrinsic dimensions.

¹From a standard uniform random variable, one can extract two independent variables by de-interleaving its binary expansion, but this operation is highly discontinuous and pathological.

3662 Bailey and Telgarsky (2018) investigate the ability of generative networks to convert input noise
3663 distributions into other distributions. One question they study is how easy it is to create a network
3664 that outputs *more* dimensions of noise than it receives. They derive bounds showing that an increase
3665 in dimension requires a large, complicated network. For example, an approximation of the uniform
3666 distribution on the unit square using the uniform distribution on the unit interval could use an
3667 (almost) space-filling curve such as the iterated tent map. (A space-filling curve is a continuous
3668 surjection from the unit interval to the unit square or volume of higher dimension.) This function is
3669 highly nonlinear and it can be shown that neural networks must be large to approximate it well.
3670 Thus the dimensionality of input noise is essential in practice. As we discuss later in this paper, our
3671 experiments support this conclusion for the game context.

3672 Bailey and Telgarsky (2018) also show that, even when the input dimension is greater than the
3673 output dimension, increased input dimension can still sometimes improve accuracy. For example, one
3674 can approximate a univariate Gaussian distribution with a high dimensional uniform distribution by
3675 summing the inputs. This follows from the **Berry-Esseen theorem** (Berry, 1941), a refinement of
3676 the central limit theorem. This uses no nonlinearity, but simply takes advantage of the fact that
3677 projecting a hypercube onto a line results in an approximately Gaussian distribution.

3678 In our experiments, we show an example of excess dimensionality improving performance.

3679 B.3 Best response computation for continuous Colonel Blotto

3680 Approximate best responses can be computed for the continuous Colonel Blotto game without
3681 discretizing the action space (Ganzfried, 2021). By sampling K batches of actions from other players'
3682 strategies, we can obtain an approximate best response a_i for player i using a mixed-integer linear
3683 program (MILP). More precisely, let h_{jk} be the highest bid for j from other players in batch k .
3684 Then solve the following MILP:

$$\text{maximize } \sum_j v_{ij} \frac{1}{K} \sum_k z_{jk} \tag{B.13}$$

$$\text{over } a_i \in \mathbb{R}^J \tag{B.14}$$

$$z \in \{0, 1\}^{J \times K} \tag{B.15}$$

$$\text{subject to } a_{ij} \geq 0 \quad \forall j \tag{B.16}$$

$$\sum_j a_{ij} = b_i \tag{B.17}$$

$$z_{jk} = [a_{ij} \geq h_{jk}] \quad \forall j, k \tag{B.18}$$

3685 We can use the Big M method to represent the last constraint:

$$a_{ij} - h_{jk} + M(1 - z_{jk}) \geq 0 \tag{B.19}$$

3686 where $M \gg 0$, which forces z_{jk} to be 0 when $a_{ij} - h_{jk}$ is negative.

3687 Appendix C

3688 Additional experiments

3689 In this section, we include some additional experiments.

3690 C.1 Ablation for Fourier features in P2SNs

3691 In this section, we study how the number of Fourier features used by a P2SN (Section 5.4) affects
3692 performance. Examples are shown in the surrounding figures. The games are as described in Section
3693 5.4.2. The figures show that a higher number of Fourier features yields superior performance. A
3694 higher number of Fourier features allows the network to represent more levels of detail for the
3695 strategies on the player set.

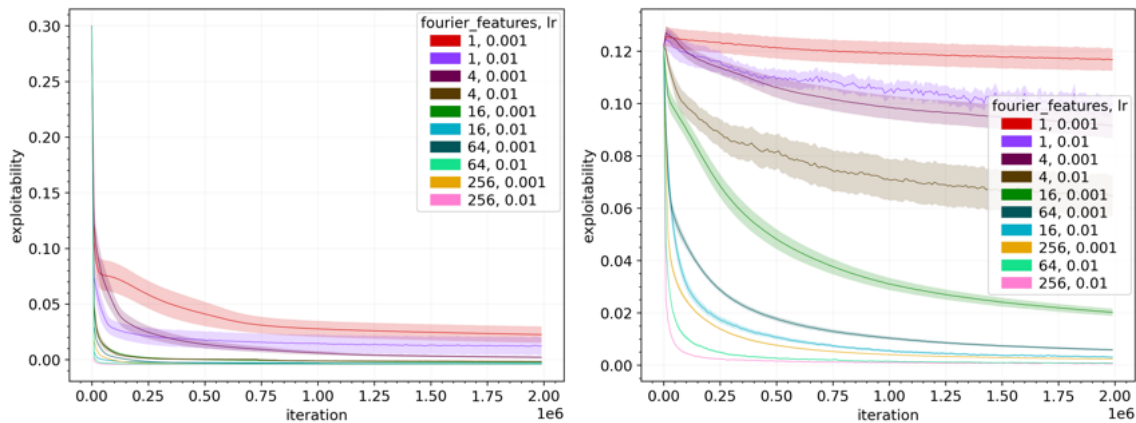


Figure C.1: Anti-coordination game. Left: 1D. Right: 2D.

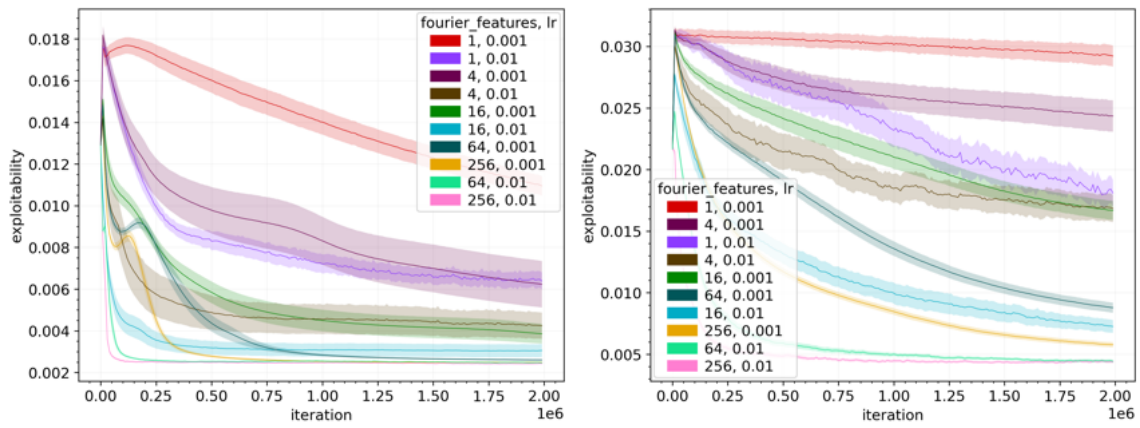


Figure C.2: Cournot game. Left: 1D. Right: 2D.

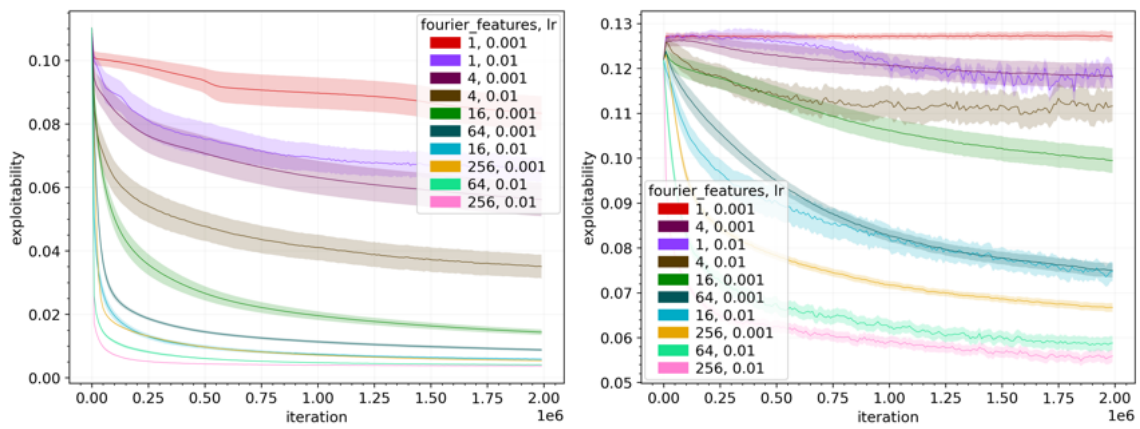


Figure C.3: Bayesian Cournot game. Left: 1D. Right: 2D.

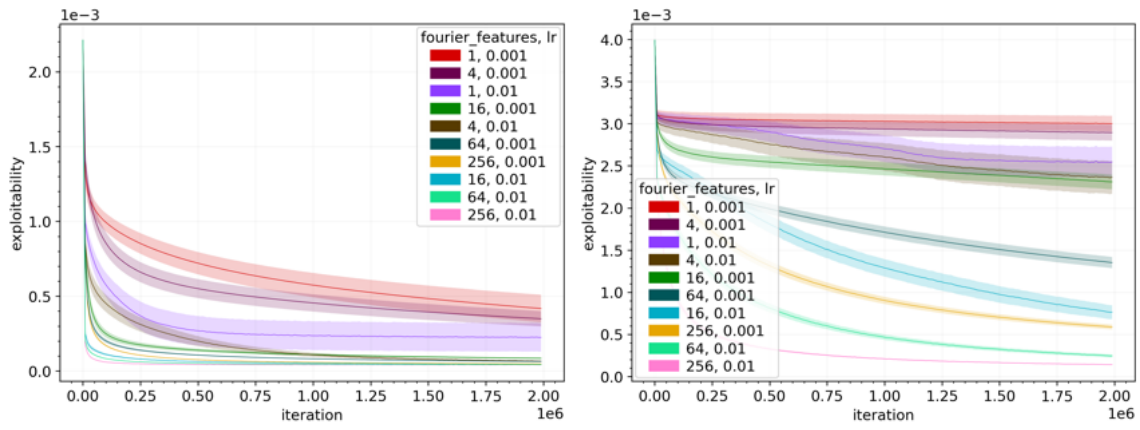


Figure C.4: Quadratic-cost Cournot game. Left: 1D. Right: 2D.

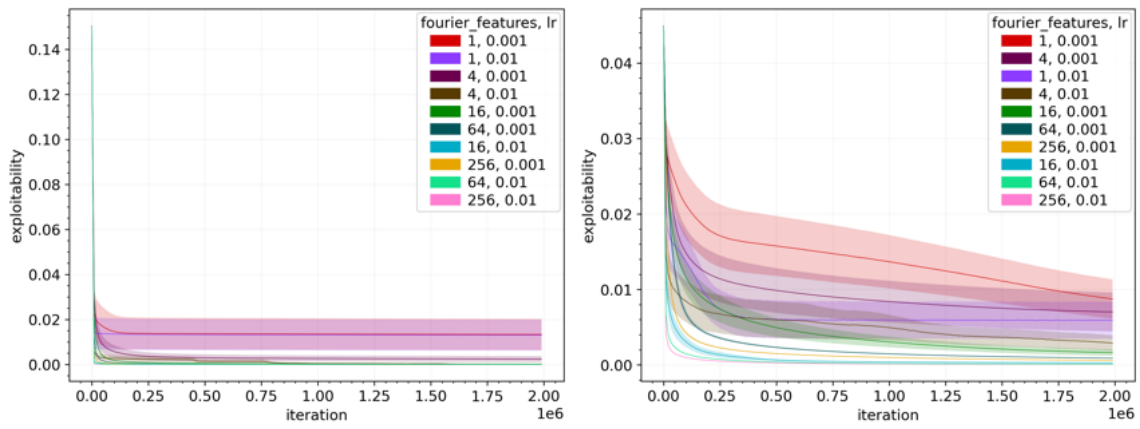


Figure C.5: Conformity game. Left: 1D. Right: 2D.

3696 C.2 Comparison to player discretization

3697 In this section, we compare P2SNs (Section 5.4) to an approach that discretizes the player space
3698 into a finite set of players and learns separate parameters for each. Specifically, given an infinite
3699 player space (such as the unit square), we use a low-discrepancy sequence (as described in Section
3700 5.4.2) to obtain a finite set of points that has good coverage of the entire space. We call these points
3701 *representatives*. To select an action for a player, we “round” it to its nearest representative (in the
3702 Euclidean metric). Thus we have a Voronoi partition of the player space. Each representative/Voronoi
3703 region has its own individual parameters that are separate from the rest.

3704 We run experiments for the anti-coordination and Cournot games described in Section 5.4.2
3705 with this approach as well as our original approach (with a width of 96 and depth of 1). Results are
3706 shown in Figures C.6–C.13. Each plot of exploitability compares our approach (labeled “p2sn”) with
3707 the discretization-based approach (labeled “discretized”). For the latter, we test various numbers of
3708 representatives (labeled “discretization_players”), thus representing different levels of resolution when
3709 discretizing the player space. We also illustrate the strategy profiles learned under the discretization
3710 approach. Visually, these approximate the strategy profiles we obtained with our approach, which
3711 were illustrated in Section 5.4.2.

3712 We observe that, across all experiments, our approach outperforms the discretization-based
3713 approach, for every tested level of resolution. We hypothesize that the reason is the following. With
3714 a finer discretization (i.e., more representatives), it takes longer to learn because *we need to learn*
3715 *more independent things*. Recall that each representative has its own individual set of parameters,
3716 and is not influenced by parameter updates to other representatives. Therefore, there are “more
3717 things that need to be learned” to obtain a good strategy profile across the entire space. Conversely,
3718 with a coarser discretization (i.e., fewer representatives), there is more parameter-sharing across
3719 players overall, but it is impossible to *represent finer details* needed for good equilibria. Thus both
3720 finer and coarser discretizations suffer from their own disadvantages. In contrast, **our method**
3721 **gets the best of both worlds!** It is able to represent finely-detailed profiles while *simultaneously*
3722 allowing players to “learn from each other” in an indirect way, due to their shared parameters.

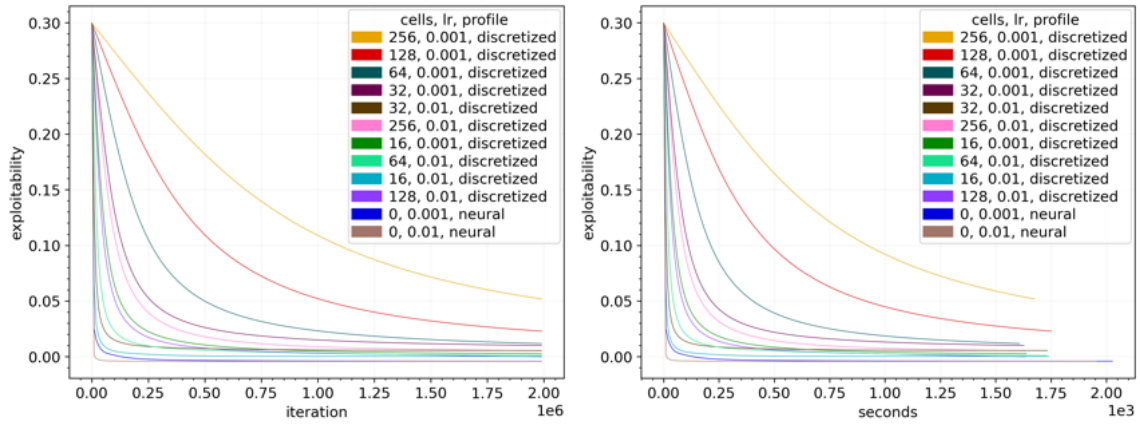


Figure C.6: 1D anti-coordination game. Exploitability. Left: Against iteration. Right: Against runtime (seconds).

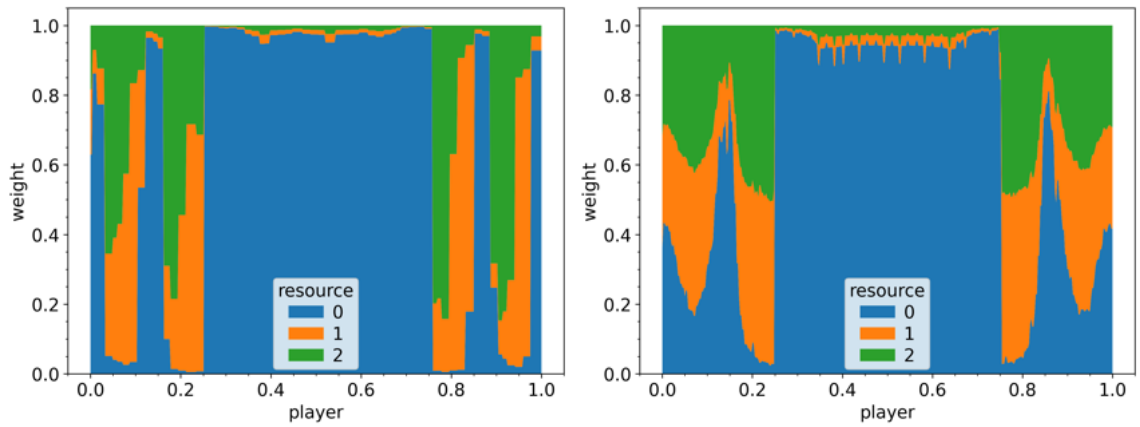


Figure C.7: 1D anti-coordination game. Learned strategy profile. Left: 64-point discretization. Right: 256-point discretization. Compare to Figure 5.39.

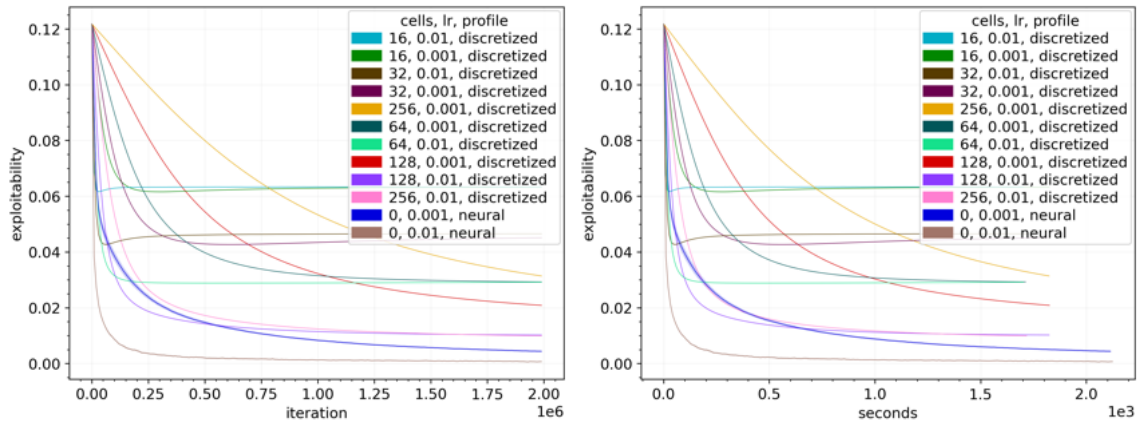


Figure C.8: 2D anti-coordination game. Exploitability. Left: Against iteration. Right: Against runtime (seconds).

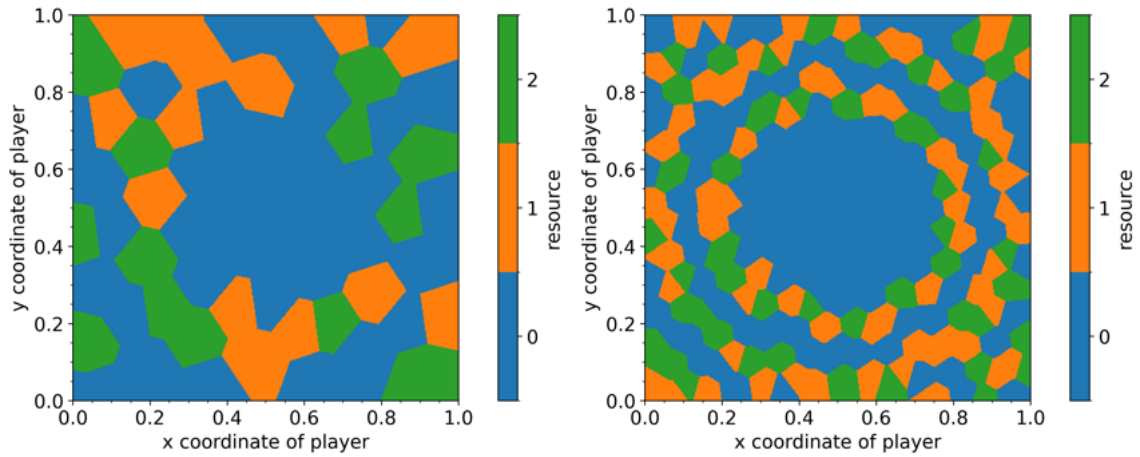


Figure C.9: 2D anti-coordination game. Learned strategy profile. Left: 64-point discretization. Right: 256-point discretization. Compare to Figure 5.40.

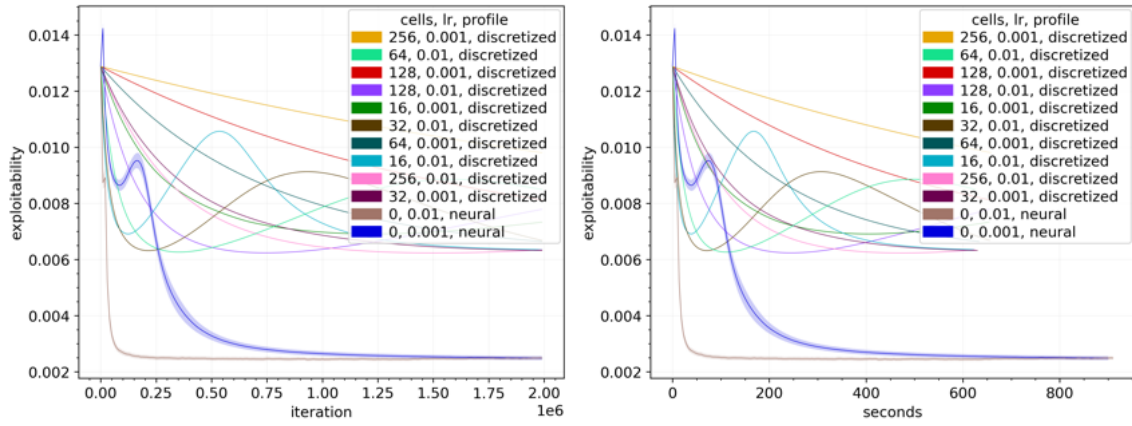


Figure C.10: 1D quadratic-cost Cournot game. Exploitability. Left: Against iteration. Right: Against runtime (seconds).

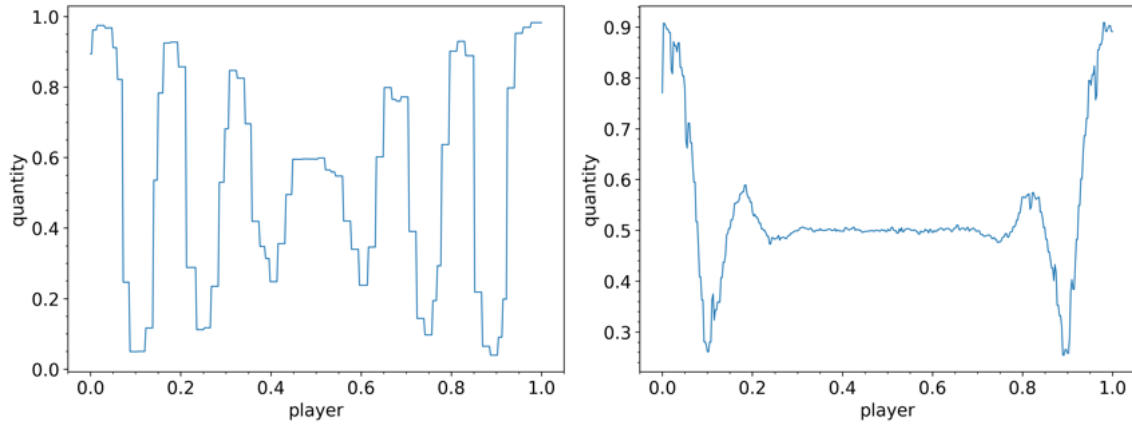


Figure C.11: 1D quadratic-cost Cournot game. Learned strategy profile. Left: 64-point discretization. Right: 256-point discretization. Compare to Figure 5.47.

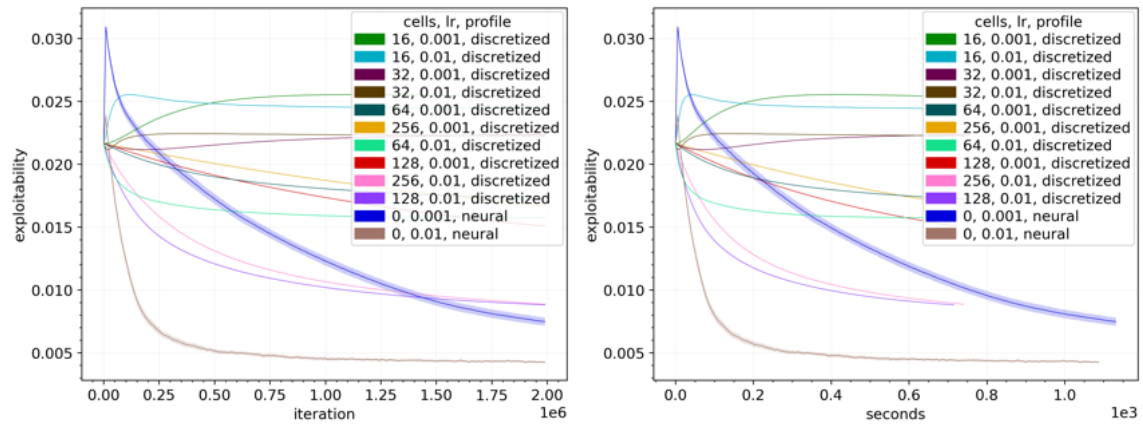


Figure C.12: 2D quadratic-cost Cournot game. Exploitability. Left: Against iteration. Right: Against runtime (seconds).

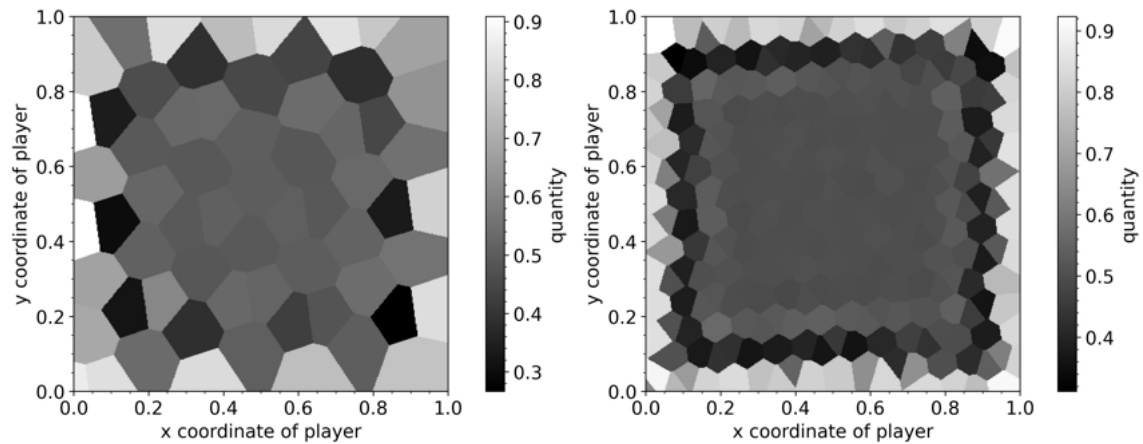


Figure C.13: 2D quadratic-cost Cournot game. Learned strategy profile. Left: 64-point discretization. Right: 256-point discretization. Compare to Figure 5.48.

3723 Appendix D

3724 Theorems

3725 Here, we review some theoretical results from the literature and our completed work.

3726 D.1 Equilibrium existence and uniqueness

3727 Here, we review theoretical results pertaining to the existence and uniqueness of NE for various
3728 kinds of games under various conditions.

3729 von Neumann (1928) proved the **minimax theorem**, showing that any two-player zero-sum game
3730 has an equilibrium in mixed strategies. This proof used Brouwer's fixed-point theorem (Brouwer,
3731 1911).

3732 Sion (1958) generalized von Neumann's minimax theorem, allowing the players' strategy sets
3733 and the utility function to have more general forms.

3734 In *Equilibrium points in n -person games*, Nash (1950) proved that any finite game with n players
3735 has an equilibrium in mixed strategies. This proof used Kakutani's fixed-point theorem (Kakutani,
3736 1941). In *Non-cooperative games*, Nash (1951) rewrote this proof to use Brouwer's fixed-point
3737 theorem.

3738 D.1.1 Infinite-action games

3739 Let (I, S, u) be an NFG, as defined in Section 3.1.

3740 Glicksberg (1952) proved that if every S_i is nonempty and compact, and $u(s)_i$ is continuous in s ,
3741 a mixed-strategy NE exists.

3742 Fudenberg and Tirole (1991, p. 34) proved that if every S_i is nonempty, compact, and convex,
3743 and $u(s)_i$ is continuous in s and quasi-concave in s_i , a pure-strategy NE exists.

3744 Dasgupta and Maskin (1986) proved that if every S_i is nonempty, compact, and convex, and
3745 $u(s)_i$ is upper semi-continuous and graph continuous in s and quasi-concave in s_i , a pure-strategy
3746 NE exists. They also proved that if S_i is nonempty, compact, and convex, and $u(s)_i$ is bounded
3747 and continuous in s except on a particular subset defined by technical conditions, weakly lower
3748 semi-continuous in s_i , and $\sum_{i \in I} u(s)_i$ is upper semi-continuous in s , a mixed-strategy NE exists.

3749 Rosen (1965) proved the uniqueness of a pure-strategy NE for continuous-action games under
3750 diagonal strict concavity assumptions.

3751 D.1.2 Infinite-player games (existence)

3752 Schmeidler (1973) generalized Nash’s theorem (Nash, 1951) on equilibrium existence to the case of
3753 a continuum of players endowed with a non-atomic measure. Khan (1985) generalized this result
3754 to non-atomic games in which each player’s strategy set is a weakly compact, convex subset of a
3755 separable Banach space whose dual has the Radon–Nikodym property. Khan (1986) generalized this
3756 result to non-atomic games with strategy sets in a Banach space.

3757 Khan and Papageorgiou (1987a) proved the existence of a Cournot–NE in generalized qualitative
3758 games with a continuum of players, under certain conditions. Similarly, Khan and Papageorgiou
3759 (1987b) proved the existence of a Cournot–NE in games with an atomless measure space of players,
3760 each with unordered preferences and strategy sets in a separable Banach space.

3761 Kim, Prikry, and Yannelis (1989) proved the existence of an NE for an abstract economy with a
3762 measure space of agents, infinite-dimensional strategy space, and agent preferences that need not be
3763 ordered.

3764 Rath (1992) gave a simple proof of the existence of pure-strategy NE in games with a continuum
3765 of players when a player’s payoff depends only on its own action and the average action of others.
3766 This was extended to the case where the action set of each player is a compact subset of \mathbb{R}^n . Khan
3767 and Sun (2002) surveyed games with many players, reporting results on the existence of pure-strategy
3768 NE in games with an atomless continuum of players, each with an action set that is not necessarily
3769 finite.

3770 Wiszniewska-Matyszek (2014) studied open- and closed-loop NE in games with a continuum of
3771 players, and both private and global state variables, proving an equivalence theorem between these
3772 classes of equilibria.

3773 Yang and Song (2022) proposed an approach for proving the existence of a pure-strategy Nash
3774 equilibrium with infinitely many players as a consequence of the finite-player case.

3775 D.1.3 Infinite-player games (uniqueness)

3776 Milchtaich (1996) proved generic uniqueness of pure-strategy NE, and uniqueness of the equilibrium
3777 outcome, for a class of non-atomic games where a player’s payoff depends on, and strictly decreases
3778 with, the measure of the set of players playing the same (pure) strategy it is playing. They also
3779 proved generic uniqueness of the Cournot–NE distribution, corresponding to a description of a game
3780 in terms of distribution of player types.

3781 A crowding game is a game in which the payoff of each player depends only on the player’s action
3782 and the size of the set of players choosing that particular action: the larger the set, the smaller the
3783 payoff. Milchtaich (2000) proved that a large crowding game generically has just one equilibrium,
3784 and the equilibrium payoffs in such a game are always unique. Moreover, the sets of equilibria of the
3785 m -replicas of a finite crowding game generically converge to a singleton as m tends to infinity.

3786 Milchtaich (2005) proved topological conditions for uniqueness of equilibrium in physical networks
3787 with a large number of users (e.g., transportation, communication, and computer networks).

3788 Caines, Huang, and Malhamé (2018) surveyed mean-field game theory. They presented its main
3789 results, including the existence and uniqueness of infinite-population NE, their approximating
3790 finite-population ε -NE, and the associated BR strategies.

3791 **Implications of equilibrium multiplicity for our method.** As previously stated, many games
3792 have a unique NE. However, some games could have multiple NE. NE refinements and equilibrium
3793 selection are outside the scope of this project, but could be an interesting question for future research.

3794 Even without such future study, the presented techniques can be very useful. In many-player games,
3795 game-theoretic approaches have been successful and often the most successful. For example, Pluribus
3796 (Brown and Sandholm, 2019) for multi-player no-limit Texas Hold'em is the only AI that has reached
3797 superhuman level in any large game beyond two-player zero-sum games. It is completely based on
3798 game-theoretic principles¹, and reached superhuman level without any training data. Many others
3799 had been trying to reach superhuman level on that exact problem for 70 years with rule-based
3800 approaches, supervised learning, reinforcement learning, etc. Only the game-theoretic approach
3801 succeeded. This is despite the game almost surely having a very large or infinite number of equilibria.

3802 D.2 Exploitability

3803 Below, $\text{convex}(x)$ and $\text{concave}(x)$ denote unspecified convex and concave functions of x , respectively.
3804 The calculation tracks convexity class through each operation. This is analogous to Big O notation
3805 (Bachmann, 1894; Landau, 1909; Knuth, 1976).

3806 **Definition 114.** A **regular game** is an NFG (I, S, u) such that $u(s)_i$ is convex in s_{-i} .

3807 **Lemma 1.** For a concave-welfare regular game, the NI is convex in the first argument.

Proof.

$$\text{NI}(x, y) = \int_{i \sim \mu} u(x[i \mapsto y_i])_i - u(x)_i \tag{D.1}$$

$$= \int_{i \sim \mu} u(x[i \mapsto y_i])_i - \int_{i \sim \mu} u(x)_i \tag{D.2}$$

$$= \int_{i \sim \mu} u(x[i \mapsto y_i])_i - \text{Wel}(x) \tag{D.3}$$

$$= \int_{i \sim \mu} u(x[i \mapsto y_i])_i - \text{concave}(x) \tag{D.4}$$

$$= \int_{i \sim \mu} u(x[i \mapsto y_i])_i + \text{convex}(x) \tag{D.5}$$

$$= \left(\int_{i \sim \mu} \text{convex}(x) \right) + \text{convex}(x) \tag{D.6}$$

$$= \text{convex}(x) + \text{convex}(x) \tag{D.7}$$

$$= \text{convex}(x) \tag{D.8}$$

3808 □

3809 **Corollary 1.** A concave-welfare regular game has convex exploitability.

Proof.

$$\text{Expl}(x) = \sup_{y \in \text{NLS}} \text{NI}(x, y) \tag{D.9}$$

¹Pluribus approximated an even weaker game-theoretic solution concept than NE, namely CCE.

$$= \sup_{y \in \Pi S} \text{convex}(x) \tag{D.10}$$

$$= \text{convex}(x) \tag{D.11}$$

3810

□

3811 **Definition 115.** A **polymatrix game** is a game with a utility function of the form

$$u(x)_i = \sum_{j \neq i} x_i^\top A_{ij} x_j \tag{D.12}$$

3812 where A_{ij} are matrices. These are graphical games in which each node corresponds to a player
 3813 and each edge corresponds to a two-player bimatrix game between its endpoints. Each player chooses
 3814 a single strategy for all of its bimatrix games and receives the sum of the resulting payoffs. In a
 3815 *constant-sum* polymatrix game, the sum of utilities across all players is constant. In a *pairwise*
 3816 *constant-sum* polymatrix game, the pairwise games are constant-sum.

3817 Cai and Daskalakis (2011) motivate the class of constant-sum polymatrix games as follows:

3818 Intuitively, these games can be used to model a broad class of competitive environments
 3819 where there is a constant amount of wealth (resources) to be split among the players of
 3820 the game, with no in-flow or out-flow of wealth that may change the total sum of players'
 3821 wealth in an outcome of the game

3822 As an example, they give a game that takes place in the wild west. A set of gold miners
 3823 need to transport gold by splitting it into wagons that traverse different paths. Each of
 3824 these paths might be controlled by thieves that could seize it.

3825 A simple example of this situation is the following game taking place in the wild west.
 3826 A set of gold miners in the west coast need to transport gold to the east coast using
 3827 wagons. Every miner can split her gold into a set of available wagons in whatever way
 3828 she wants (or even randomize among partitions). Every wagon uses a specific path to go
 3829 through the Rocky mountains.

3830 Unfortunately each of the available paths is controlled by a group of thieves. A group of
 3831 thieves may control several of these paths and if they happen to wait on the path used
 3832 by a particular wagon they can ambush the wagon and steal the gold being carried. On
 3833 the other hand, if they wait on a particular path they will miss on the opportunity to
 3834 ambush the wagons going through the other paths in their realm as all wagons will cross
 3835 simultaneously.

3836 The utility of each miner in this game is the amount of her shipped gold that reaches
 3837 her destination in the east coast, while the utility of each group of thieves is the total
 3838 amount of gold they steal. Clearly, the total utility of all players in the wild west game
 3839 is constant in every outcome of the game (it equals the total amount of gold shipped by
 3840 the miners), but the pairwise interaction between every miner and group of thieves is
 3841 not. In other words, the constant-sum property is a *global* rather than a *local* property
 3842 of this game.

3843 Further applications and a discussion of several special cases of these games can be found in
 3844 Bergman and Fokin (1998).

3845 Cai and Daskalakis (2011) prove a generalization of von Neumann’s minmax theorem to constant-
3846 sum polymatrix games. Their theorem implies convexity of NE, polynomial-time tractability, and
3847 convergence of no-regret learning algorithms to NE.

3848 Cai, Candogan, et al. (2016) show that, in such games, NE can be found efficiently via linear
3849 programming. They also show that the set of CCE collapses to the set of NE.

3850 We prove that such games have convex exploitability.

3851 **Theorem 1.** A constant-sum polymatrix game has convex exploitability.

3852 *Proof.* For $j \neq i$, $x_i^\top A_{ij} x_j$ is linear, and therefore convex, in x_j . Therefore, their sum $u(x)_i$ is convex
3853 in x_{-i} . Furthermore, the game is constant-sum. Since the game is constant-sum, the welfare is
3854 constant, hence concave. Thus, by Corollary 1, the exploitability is convex. \square

3855 **Corollary 2.** A two-player constant-sum matrix game has convex exploitability.

3856 Separately, we have the following result:

3857 **Theorem 2.** A two-player constant-sum concave-convex game has convex exploitability.

3858 *Proof.* The exploitability reduces to the *duality gap*, as described in Grnarova, Kilcher, et al. (2021):

$$\text{Expl}(x) = \sup_{y_1 \in S_1} u(y_1, x_2)_1 - u(x)_1 + \sup_{y_2 \in S_2} u(x_1, y_2)_2 - u(x)_2 \quad (\text{D.13})$$

$$= \sup_{y_1 \in S_1} u(y_1, x_2)_1 - \inf_{y_2 \in S_2} u(x_1, y_2)_2 - u(x)_1 - u(x)_2 \quad (\text{D.14})$$

$$= \sup_{y_1 \in S_1} u(y_1, x_2)_1 - \inf_{y_2 \in S_2} u(x_1, y_2)_1 + \text{constant}(x) \quad (\text{D.15})$$

$$= \sup_{y_1 \in S_1} \text{convex}(x) - \inf_{y_2 \in S_2} \text{concave}(x) + \text{constant}(x) \quad (\text{D.16})$$

$$= \text{convex}(x) - \text{concave}(x) + \text{constant}(x) \quad (\text{D.17})$$

$$= \text{convex}(x) \quad (\text{D.18})$$

3859

\square

3860 D.3 Subgradient descent

3861 Our techniques seek to minimize (an approximation of) exploitability by performing subgradient
3862 descent. This raises the question of when such a process is able to attain a global minimum in the
3863 first place.

3864 Kiwiel (2004) analyze the convergence of **approximate subgradient methods** for convex
3865 optimization, and prove the following theorems. Let $\mathcal{S} \subseteq \mathbb{R}^n$ be a nonempty closed convex set,
3866 $f : \mathcal{S} \rightarrow \mathbb{R}$ be a closed proper convex function, and $\mathcal{S}_* \in \text{argmin } f$. Let $x_{t+1} = P_{\mathcal{S}}(x_t - \nu_t g_t)$ where $P_{\mathcal{S}}$
3867 is the projector onto \mathcal{S} ($P_{\mathcal{S}}(x) \in \text{argmin}_{y \in \mathcal{S}} \|x - y\|$), $\nu_t \geq 0$ is a stepsize, $\varepsilon_t \geq 0$ is an error tolerance,
3868 and $g_t \in \partial_{\varepsilon_t} f(x_t)$ is an ε_t -approximate subgradient of f at x_t , that is, $f(x) \geq f(x_t) + \langle g_t, x - x_t \rangle - \varepsilon_t$
3869 for all x .

3870 **Theorem 3.** (Kiwiel, 2004, Theorem 3.4) Suppose $\mathcal{S}_* \neq \emptyset$, $\sum_{t \in \mathbb{N}} \nu_t = \infty$, and $\sum_{t \in \mathbb{N}} \nu_t (\frac{1}{2} \|g_t\|^2 \nu_t +$
3871 $\varepsilon_t) < \infty$. Then $\{x_t\}_{t \in \mathbb{N}}$ converges to some $x_\infty \in \mathcal{S}_*$.

3872 **Theorem 4.** (Kiwiel, 2004, Theorem 3.6) Suppose $\mathcal{S}_* \neq \emptyset$, $\sum_{t \in \mathbb{N}} \nu_t = \infty$, $\sum_{t \in \mathbb{N}} \nu_t^2 < \infty$,
3873 $\sum_{t \in \mathbb{N}} \nu_t \varepsilon_t < \infty$, and the subgradients do not grow too fast: $\exists c < \infty, \forall t \in \mathbb{N}, \|g_t\|^2 \leq c(1 + \|x_t\|^2)$
3874 (e.g., they are bounded). Then $\{x_t\}_{t \in \mathbb{N}}$ converges to some $x_\infty \in \mathcal{S}_*$.

3875 Convergence results are also known for subgradient methods on *quasi*-convex functions, which is
3876 a more general class of functions. Some of these are described in Hu, Yang, and Sim (2015).

3877 In our paper, we seek to approximately minimize the exploitability function $\Phi \in \mathcal{S} \rightarrow \mathbb{R}$,
3878 $\Phi(x) = \sup_{y \in \mathcal{S}} \phi(x, y)$, where ϕ is the Nikaido–Isoda function and \mathcal{S} is the set of possible strategy
3879 profiles. Specifically, we use subgradients of $\tilde{\Phi}_t(x) = \phi(x, \tilde{y}_t)$, where \tilde{y}_t is the response profile output
3880 by the learned **best-response ensembles** (BRE) or **best-response function** (BRF) at t . Thus
3881 $\tilde{\Phi}_t(x) \geq \tilde{\Phi}_t(x_t) + \langle g_t, x - x_t \rangle$. Unconditionally, $\Phi \geq \tilde{\Phi}_t$ (since the former maximizes over all possible
3882 y). Thus $\Phi(x) \geq \tilde{\Phi}_t(x_t) + \langle g_t, x - x_t \rangle$.

3883 Now, suppose we can guarantee that $\tilde{\Phi}_t(x_t) \geq \Phi(x_t) - \varepsilon_t$ for an error tolerance $\varepsilon_t \geq 0$; that is,
3884 the responses output by the BRE/BRF at t do not perform *too* badly (in the limit) compared to
3885 the true best responses. Then $\Phi(x) \geq \Phi(x_t) - \varepsilon_t + \langle g_t, x - x_t \rangle$.

3886 Therefore, when the assumptions of the above theorems hold, the sequence of iterates $\{x_t\}_{t \in \mathbb{N}}$
3887 converges to a global minimizer of the exploitability function, which is an NE, if an NE exists at all.

3888 D.4 Theoretical results pertaining to P2SN

3889 In this section, we include some theoretical analysis of the method we proposed in Martin and
3890 Sandholm (2025e).

3891 D.4.1 Progress guarantee

3892 Let

$$g = \left\{ \frac{du(s_\theta)_i}{d(s_\theta)_i} \right\}_{i \in P} \quad (\text{D.19})$$

3893 be the true simultaneous gradient in the (potentially infinite-dimensional) **function space** of
3894 strategy profiles. Let

$$v = \left[\int_{i \sim \mu} \frac{d}{d\theta} u(s_\phi[i \mapsto (s_\theta)_i])_i \right]_{\phi=\theta} \quad (\text{D.20})$$

$$= \int_{i \sim \mu} \frac{d(s_\theta)_i}{d\theta} \frac{du(s_\theta)_i}{d(s_\theta)_i} \quad (\text{D.21})$$

3895 be the SPSG. The change in strategy profile induced by the SPSG is

$$\tilde{g} = Jv \quad (\text{D.22})$$

3896 where $J = \frac{ds_\theta}{d\theta}$ is the network Jacobian. Furthermore,

$$v = \int_{i \sim \mu} J_i g_i \quad (\text{D.23})$$

$$= J^\top g \quad (\text{D.24})$$

3897 Therefore,

$$g \cdot \tilde{g} = g \cdot (Jv) \quad \text{Eq. D.22} \quad (\text{D.25})$$

$$= (J^\top g) \cdot v \quad (\text{D.26})$$

$$= v \cdot v \quad \text{Eq. D.24} \quad (\text{D.27})$$

$$= \|v\|^2 \quad (\text{D.28})$$

$$\geq 0 \quad (\text{D.29})$$

3898 That is, the true simultaneous gradient and SPSG-induced gradient are compatible. Furthermore,

$$v \neq 0 \leftrightarrow g \notin \text{kernel}(J^\top) \leftrightarrow g \notin \text{im}(J) \quad (\text{D.30})$$

3899 Therefore,

$$g \cdot \tilde{g} > 0 \leftrightarrow g \notin \text{kernel}(J^\top) \leftrightarrow g \notin \text{im}(J) \quad (\text{D.31})$$

3900 That is, if the simultaneous gradient has a nonzero projection onto the tangent space of the network—
 3901 i.e., onto the Jacobian—progress is made. Any components of the true simultaneous gradient that
 3902 are orthogonal to the tangent space cannot be represented by any parameter update. If $g \neq 0$ but
 3903 $v = 0$, then the network-representable tangent space cannot implement any average improvement
 3904 direction. In practice, this means that the network must be expressive enough to move in directions
 3905 that improve players' utilities on average.

3906 D.4.2 Special cases

3907 In this section, we develop stronger conclusions under various additional assumptions about the
 3908 setting.

3909 **Disjoint parameters.** Suppose $(s_\theta)_i$ depends only on θ_i . Then

$$v(\theta) = \int_{i \sim \mu} \left. \frac{du(s_\phi[i \mapsto (s_\theta)_i])_i}{d\theta} \right|_{\phi=\theta} \quad (\text{D.32})$$

$$= \int_{i \sim \mu} \left. \frac{du(s_\phi[i \mapsto (s_\theta)_i])_i}{d\theta_i} \right|_{\phi=\theta} \frac{d\theta_i}{d\theta} \quad (\text{D.33})$$

$$= \int_{i \sim \mu} \frac{du(s_\theta)_i}{d\theta_i} \mathbf{e}_i \quad (\text{D.34})$$

3910 Thus SPSG reduces to the usual simultaneous gradient.

3911 **Potential game.** Suppose the game is a potential game. Let ϕ be its potential function. Then

$$u(s[i \mapsto a])_i - u(s[i \mapsto b])_i = \phi(s[i \mapsto a]) - \phi(s[i \mapsto b]) \quad (\text{D.35})$$

3912 Thus

$$\frac{du(s)_i}{ds_i} = \frac{d\phi(s)}{ds_i} \quad (\text{D.36})$$

3913 Let L be the potential composed with the strategy network.

$$L(\theta) = \phi(s_\theta) \quad (\text{D.37})$$

3914 On the one hand,

$$\frac{dL(\theta)}{d\theta} = \frac{d\phi(s_\theta)}{d\theta} \quad (\text{D.38})$$

$$= \int_{i \sim \mu} \frac{d\phi(s_\theta)}{d(s_\theta)_i} \frac{d(s_\theta)_i}{d\theta} \quad (\text{D.39})$$

$$= \int_{i \sim \mu} \frac{du(s_\theta)_i}{d(s_\theta)_i} \frac{d(s_\theta)_i}{d\theta} \quad \text{Eq. D.36} \quad (\text{D.40})$$

3915 On the other hand,

$$v(\theta) = \int_{i \sim \mu} \left. \frac{du(s_\theta[i \mapsto (s_\theta)_i])_i}{d\theta} \right|_{\phi=\theta} \quad (\text{D.41})$$

$$= \int_{i \sim \mu} \left. \frac{du(s_\theta[i \mapsto (s_\theta)_i])_i}{d(s_\theta)_i} \right|_{\phi=\theta} \frac{d(s_\theta)_i}{d\theta} \quad (\text{D.42})$$

$$= \int_{i \sim \mu} \frac{du(s_\theta)_i}{d(s_\theta)_i} \frac{d(s_\theta)_i}{d\theta} \quad (\text{D.43})$$

3916 Therefore

$$v(\theta) = \frac{dL(\theta)}{d\theta} \quad (\text{D.44})$$

3917 That is, SPSPG is equivalent to stochastic gradient ascent on the network-transformed potential.
 3918 For example, if ϕ is convex and $\theta \mapsto s_\theta$ is affine in θ , L is also convex. All theoretical guarantees
 3919 for stochastic gradient ascent carry over to this setting, given appropriate bounded-variance and
 3920 integrability conditions on the integrand of $v(\theta)$.

3921 **Concave utilities.** Suppose that each player's utility is concave and differentiable in its own
 3922 strategy. Let $\ell = -u$ be the loss function. (Here, we work with losses because it is more common in
 3923 the optimization literature.) Then each player's loss is convex and differentiable in its own strategy.
 3924 Suppose that $(s_\theta)_i = A_i\theta + b_i$ is affine. Thus $\nabla_{s_i}^2 \ell(s)_i \succeq 0$ for all s .

3925 Let $\tilde{\ell}(\theta) = \ell(r[i \mapsto (s_\theta)_i])_i$, where r is the strategy profile that player i is deviating from. Then

$$\nabla_\theta \tilde{\ell}(\theta) = (\nabla_\theta (s_\theta)_i) (\nabla_{(s_\theta)_i} \tilde{\ell}(\theta)) \quad (\text{D.45})$$

$$= A_i^\top (\nabla_{(s_\theta)_i} \tilde{\ell}(\theta)) \quad (\text{D.46})$$

3926 and

$$\nabla_\theta^2 \tilde{\ell}(\theta) = A_i^\top (\nabla_{(s_\theta)_i}^2 \tilde{\ell}(\theta)) A_i \quad (\text{D.47})$$

$$= A_i^\top (\nabla_{(s_\theta)_i}^2 \ell(s_\theta)) A_i \quad (\text{D.48})$$

3927 By hypothesis $\nabla_{s_i}^2 \ell(s)_i \succeq 0$, thus $A_i^\top (\nabla_{s_i}^2 \ell(s)_i) A_i \succeq 0$. Therefore, the symmetric part of $\nabla_\theta^2 \tilde{\ell}(\theta)$
 3928 is PSD. Integrating over i preserves this, so the symmetric part of $\nabla_\theta v(\theta)$ is PSD. Since $v(\theta)$ is
 3929 continuously differentiable and its symmetric Jacobian is PSD, it is a monotone operator. Thus, the
 3930 theoretical convergence guarantees for stochastic monotone operators apply. Loizou et al. (2021)
 3931 gives an example of a general convergence guarantee for such operators.