**FACULTY**
**OF MATHEMATICS**
**AND PHYSICS**
**Charles University**

**DOCTORAL THESIS**

Martin Schmid

# Search in Imperfect Information Games

Department of Applied Mathematics

| | |
|---|---|
| Supervisor of the doctoral thesis: | Milan Hladík |
| Advisor of the doctoral thesis: | Michael Bowling |
| Study programme: | Computer Science |
| Study branch: | Discrete Models and Algorithms |

Prague 2021

I declare that I carried out this doctoral thesis independently, and only with the cited sources, literature and other professional sources.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In ........ date ............          signature of the author

Title: Search in Imperfect Information Games

Author: Martin Schmid

Department: Department of Applied Mathematics

Supervisor: Milan Hladík, Department of Applied Mathematics

Advisor: Michael Bowling, Department of Computing Science, University of Alberta

Abstract: From the very dawn of the field, search with value functions was a fundamental concept of computer games research. Turing's chess algorithm from 1950 was able to think two moves ahead (Copeland, 2004), and Shannon's work on chess from 1950 includes an extensive section on evaluation functions to be used within a search (Shannon, 1950). Samuel's checkers program from 1959 already combines search and value functions that are learned through self-play and bootstrapping (Samuel, 1959). TD-Gammon improves upon those ideas and uses neural networks to learn those complex value functions — only to be again used within search (Tesauro, 1995). The combination of decision-time search and value functions has been present in the remarkable milestones where computers bested their human counterparts in long standing challenging games — DeepBlue for Chess (Campbell et al., 2002) and AlphaGo for Go (Silver et al., 2016).

Until recently, this powerful framework of search aided with (learned) value functions has been limited to perfect information games. As many interesting problems do not provide the agent perfect information of the environment, this was an unfortunate limitation. This thesis introduces the reader to sound search for imperfect information games.

Keywords: game theory, search, imperfect information, games, DeepStack

# Acknowledgments

There are many people who I am deeply grateful to and for.

My supervisor Milan Hladik allowed for unprecedented level of freedom in my research. He trusted me to pursue ideas and goals that I found interesting and relevant, always believing me to achieve the challenging dreams of mine. He also included me in many unforgettable seminars at the society for the optimization of the consumption, experience which one can only enjoy.

I have met my advisor Michael Bowling during my very first conference and my very first travel across the ocean. At the same conference two years later is where we discussed the exciting plan of me and Matej visiting the Edmonton office to join his team. I vividly remember me pacing back and forth at the hotel room, pondering the idea of moving to Canada to pursue the dream of DeepStack. Since then, my work under Mike has deeply molded me and made me a better person and better scientist. As his motto is "once my student, always my student", I am proud to belong to this club.

DeepStack brings me to Matej Moravcik, with whom I have spent sleepless nights — full of coffee, code and math (the most enjoyable combination really!), working for almost a year to finish the project. Never will I forget the thrill of working on this amazing break-through, the ups and downs, the excitement of pushing the boundaries and the scientific joy the live is worth living for.

My scientific mind and heart was from the very early age forged by my mother. She is the archetypal absent-minded professor who toughed me to be curious, think scientifically and work hard. To never take yourself too seriously, to ignore the mundane standards of the masses but to do the things that excite me and make me happy.

Only in silence the word, only in dark the light, only in dying life: bright the hawk's flight on the empty sky.

Ursula K. Le Guin, A Wizard of Earthsea

# Contents

# Contributions

This thesis aims to tell a coherent story of search in imperfect information games. My main contributions in this line of work are described below. Each contribution has been published in a major conference or a journal (see given citations and the full list of author's publications at the end of the thesis in Section 21.7.5).

**Bounding Support Sizes in Games**  In order to play optimally, imperfect information games often require stochastic policies. This is in contrast to perfect information games, where an optimal deterministic strategy always exist. We have proven an intriguing connection between a players' level of uncertainty in state of the game and the support size (number of actions with non-zero probability) (Schmid et al., 2014). Section 14.5.1 includes a brief overview of the result as most of the contribution dates to an older version of the paper, published prior to the beginning of my doctoral study (Schmid and Moravcik, 2013).

**Factored Observation Games Formalism**  The work on new sound search methods for imperfect information games revealed that previous formalisms are ill-suited for this task. We have introduced a new formalism particularly suited for the needs of modern search algorithms (Kovařík et al., 2019). We have also proven connections to the prior formalism of extensive form games. Finally, we formulated counterfactual regret minimization (Section 17.5) and sequence form linear optimisation (Section 16.3) in this formalism — popular algorithms from n extensive form games. This formalism is also adopted for most of this thesis. See also Section 13.3.

**Variance reduced Monte Carlo Counterfactual Regret Minimization**  Monte Carlo CFR (MCCFR) is a family of game tree sampling algorithms for imperfect information games that has become a key-building block of many recent methods. As the algorithm does not traverse the entire game tree but rather samples trajectories, the corresponding sampled values are inherently noisy. While MCCFR still provides strong (probabilistic) convergence bounds, the variance introduced by sampling can greatly slow down the convergence speed. We have introduced VR-MCCFR that decreases the variance and results in orders of magnitude faster convergence (Schmid et al., 2019; Davis et al., 2020). For details see Section 17.8.

**Refining Sub-games in Large Imperfect Information Games**  Prior to the dawn of search, state of the art methods were based on the abstraction framework which simply solved a smaller, abstracted game (Section 18.1). We have introduced sound and safe refinement of sub-games, that allows to on-line improve the strategy near the end of the game when the corresponding sub-game becomes (more) tractable (Moravcik et al., 2016). More details in Section 20.4.

**Formalising $\epsilon$-Sound Search**  The key solution concept of $\epsilon$-Nash equilibrium is defined for fixed, offline strategies. We have generalized the concept for online settings, allowing one to analyze the worst-case behavior of search algorithms

(Šustr et al., 2020, 2021). Our consistency hierarchy then revealed further discrepancies between perfect and imperfect information games. Relevant Chapters are 7 and 19.

**DeepStack** The culmination of the presented thesis is the sound search technique of continual re-solving (Section 20.7). Continual re-solving with value functions represented by deep neural networks is the algorithm behind the DeepStack agent (`https://www.deepstack.ai`). DeepStack (Chapter 21) became the first agent to beat professional human players in the game of no-limit Texas hold'em poker, combining search and learning and constitutes a major break-through for search in imperfect information games (Moravčík et al., 2017).

**AIVAT** As poker is an inherently high variance game, many matches are often required to get a statistically significantly result. While this is possible for computer agents where computer poker competitions often required millions of data-points during evaluation, it is expensive and even more problematic when human play is involved. We have developed AIVAT (Burch et al., 2018), a provably unbiased technique for reducing the variance during evaluation. Note that this technique is not poker specific and can be used to evaluate any game. AIVAT and prior variance reduction techniques is described in Section 21.6

# 1. Introduction

Games have long served as benchmarks and marked milestones of progress in artificial intelligence (AI), serving as "fruit flies" of AI research (McCarthy, 1990). The same way lab mice allow the researchers to speed up a development of new human drugs, games allow us to design algorithms for challenging real world environments. Unlike mice, games allow for an easy natural progression, as different games bring varying challenges. Games can be single player or multiplayer, perfect or imperfect information, simultaneous moves or sequential, discrete or continuous actions, etc.

Last but not least, games are fun. Games are fun to play, and even more fun to do research in. "Don't worry about the overall importance of the problem; work on it if it looks interesting. I think there is a sufficient correlation between interest and importance" (David Blackwell).

## 1.1   Learning

As we hope that games help us to design algorithms for a wide class of problems, the more general the algorithms are, the better. If one had one algorithm for chess, one for poker an another one for security games, it would be harder to imagine how these techniques seamlessly generalize to other problems.

A powerful idea that is appealing to incorporate into our approaches is learning. Rather than hard-coding specific game algorithm, we can design an algorithm that learns to play the game by repeatedly interacting with the game and improving its policy.

## 1.2   Games in Question

This thesis is limited to a particular subset of games. The class of games we will be concerned with are strictly adversarial (zero sum) two player games. Many popular games fit these criteria (e.g. chess, poker, go etc.). Search for this relatively large class of games is still a step forward from perfect information games, as imperfect information games generalize perfect information games.

Why this particular class? First, the optimal policies provide strong and appealing properties. Second, these optimal policies are computable in polynomial time (in the size of the game). Third, it is not clear what policies would be desirable in the larger class of multiplayer games. The common solution concepts lose the appealing guarantees and thus their appeal. Overall, the theoretical results for multiplayer games are rather unfortunate: equilibrium policies are harder to compute and they lack the same appealing guarantees as in two-player case.

## 1.3   Optimal Policies

In single agent environments, the common concept of optimal policy is that of a policy that maximizes the return. In the case of multiple agents, the reward depends not just on the actions of our agent, but also on the actions of all the

other players in the game. As we do not get to know the policy of the opponents, the question of optimality is complicated. We still want to maximize return, but how can we do that as we do not know what the opponent will do?

### 1.3.1 Worst Case Opponent — Maximin

One option is to consider the worst case opponent. This is rather natural if we think of how one reasons in chess. Consider an overly simplistic situation where we simply need to maximize chess material. When making a chess move, would you rather make:

- A move that no matter what counter-move the opponent does, you will win a pawn.

- A move that for almost all of the counter-moves of the opponent, you will win a rook. But there is one counter-move of the opponent that loses you a queen.

Optimizing the reward against a worst-case opponent will lead to the popular maximin solution concept. This solution concept is guaranteed to exist in both perfect and imperfect information games, with some appealing properties.

### 1.3.2 Nash Equilibrium

We will also see another solution concept, Nash equilibrium. Unlike the worst case reasoning, Nash equilibrium is defined as a joint policy (strategy profile) for all the players where none of the players benefit from deviating from this profile. Crucially, we will learn that for the games in question — two player zero sum perfect and imperfect information games — the concepts of minmax and Nash equilibria collapse. We thus often use the term optimal policy to refer to both solution concepts interchangeably.

### 1.3.3 The Appeal of Optimal Policies

Arguably the most appealing property of optimal policies in the context of two-player zero-sum games is that if we get to play both positions (e.g. white and black in chess, small and big blind in poker), we can not lose on expectation against any opponent. This allows us to focus on algorithms that scale to large problems and produce provably near-optimal strategies, rather than dealing with the distributions of the opponent we could be facing.

## 1.4 Offline Algorithms

Offline algorithms compute the optimal policy prior to the actual game play, in other words, solve the game. The result of the computation then describes what policy to follow in each state of the game, and is then simply retrieved during the game-play.

### 1.4.1 Tabular vs Implicit Representation

An important distinction has to be made between tabular methods and methods that store the strategy in some implicit way. A natural representation of the computed strategy is a tabular one — that is, storing explicitly how to act in every state $s \in \mathcal{S}$. While such a representation might be necessary for an optimal strategy, it can be intractable for large games.

It is worth mentioning that in some cases, it is possible to represent provably optimal policies without the need to the store strategy for each state of the game. Consider the idea of alpha-beta pruning that allows one to prune parts of the game tree in a sound way (Knuth and Moore, 1975). The resulting policy can then be thought of as a small certificate and one can verify the validity of such a certificate in time linear in the size of the certificate (in contrast to the size of the full game). Interestingly, this idea can be generalized to Nash equilibrium certificates in imperfect information games (Zhang and Sandholm, 2020).

Rather than explicit representation of the policy, one can opt for an implicit representation, for example by a parametric function representation (a state to a strategy). If we represent the states by a single number $s \in \{1 \dots |S|\}$ and the actions in all the states are $\{a, b\}$, we can represent a policy by letting $p(a|s) \propto \frac{1}{s}$. While such representation is rather silly, it illustrates the idea. In practice, the function would be some more complex parametric function to produce $p(a|s) = f_\Theta(a, s)$, maybe represented by a deep neural network.

## 1.5 Search

Search is a particular form of an online algorithm, similar to the way humans play chess (Newell and Simon, 1964). Unlike offline algorithms, online algorithms compute the policy for the required/observed state during game-play. Given a state of a game, search algorithms reason about potential future states using a look-ahead tree. As the number of all future states can be too large, this look-ahead is often truncated and the terminal states assigned a heuristics value. This evaluation ideally would closely approximate the true value of these states under the assumption of optimality. That is, the value of a future state is the value of the sub-problem rooted in that state, given that both players play optimally in that sub-problem. Search algorithms then reason about optimal play within the look-ahead tree using these sub-problem values, resulting in an optimal policy for the state in question. Search is thus typically a combination of i) sub-problem decomposition, ii) (learning of) value functions iii) local (look-ahead) reasoning method.

### 1.5.1 Benefits

There are multiple reasons search is an appealing method.

First, we get to use compute resources during the game-play and can focus these resources on the part of the problem that is immediately relevant. Search algorithms compute a strategy for the current state to play and do not have to reason about the states irrelevant to the current situation. Purely offline

algorithms do not know what states will be visited during the play, and thus have to compute the strategy for all of them.

Second, search methods get to leverage a model of the environment. This allows it to exactly reason about the future states and the estimated values (or even exact values in the case of terminal states). While the model can be potentially learned (Schrittwieser et al., 2020), we limit this work to settings with an exact game model.

While these reason sound appealing, what matters most is empirical performance and search methods do in fact excel in practice. All the top computer agents in chess, go, shogi, Arimaa and many other games share one thing: they all use search methods. And now the same argument of performance holds for imperfect information games such as poker.

## 1.6 Imperfect Information Games



Figure 1.1: a) Perfect information chess. b) Imperfect information chess. c) Simultaneous moves chess. Chessboard image under the creative commons license[1]

This thesis deals with search in imperfect information games. The defining property is that unlike in perfect information games (e.g. chess, go), agents do not have a full (or symmetric) information about the state of the game. But why is this imperfect information property of any importance? Consider the following three chess variants (where the first one is perfect information game, while the other two are imperfect information games).

1. Standard, perfect information chess (Figure 1.1a).

2. Imperfect information variant, where the players do not see the opponent's pieces in the the shaded squares (Figure 1.1b).

3. Simultaneous moves chess, where the players make move at the same time (Figure 1.1c).

### 1.6.1 Optimal Policies

All these chess variants make no difference when it comes to the existence or guarantees of optimal policies. We will see that just like in perfect information

---

[1]AAA SVG Chessboard and chess pieces06 from ILA-boy, licensed under CC BY-SA 3.0

chess, optimal policies still exist and have all the same desirable properties. One distinction of the optimal policies in imperfect information games is that it often has to be stochastic. While there always is a deterministic optimal policy in perfect information games, this is not the case for the other variants. It turns out that the number of actions we need to randomize between is closely related to the level of uncertainty in the game (Schmid et al., 2014).

From the algorithmic perspective, the crucial distinction is that finding optimal policies is more challenging. Many offline algorithms that provably work in perfect information, fail to find optimal policies in imperfect information games, most notably many popular[2] reinforcement learning algorithms.

### 1.6.2 Search

As offline algorithms for imperfect information games are more challenging, so are the online algorithms that do search. Search consists of three fundamental components: i) local reasoning ii) sub-problem decomposition iii) value functions. As we will see, all of these components are substantially more complex in imperfect information settings. While some popular perfect information based methods (e.g. Monte Carlo tree search) can be used in imperfect information (Whitehouse, 2014), these methods are not sound and fail to produce optimal policies even in very small games. For a long time, sound search has been thought to be impossible in imperfect information settings (Frank et al., 1998; Lanctot et al., 2014).

This thesis tries to coherently introduce the reader to a history of computational game theory research leading to this important milestone: sound search in imperfect information.

## 1.7 Reinforcement Learning or Game Theory

There are two main disciplines and corresponding communities that study sequential decision-making in multi-agent environments using games as benchmarks for progress: reinforcement learning and game theory. In reinforcement learning, search and learning of value functions quickly became a common ingredient of many algorithms (Sutton et al., 1998).

On the other hand, game theory typically focused on solution concepts, complexity results and algorithms with strong theoretical guarantees - producing provably optimal policies for imperfect information games (Shoham and Leyton-Brown, 2008).

Are the techniques presented reinforcement learning or game theory techniques? It is the opinion of the authors that they are both. The techniques combine value functions, learning, search and strong theoretical convergence guarantees in both perfect and imperfect information games. Furthermore, the line between both communities is blurring and many challenging environments will require insights of both communities.

---

[2](Deep) Q-learning, SARSA, REINFORCE, PPO etc.

### 1.7.1 Dictionary

While both communities deal with optimal behavior of agent(s) in an environment and thus use the same concepts, the language even for the basic concepts often differs. This Babylonian confusion of tongues has mostly historical reasons. But if we want to finish the tower of general algorithms, we need the work to be accessible to both communities. Table 1.1 lists some of the basic concepts and we purposely use words from both reinforcement learning and game theory throughout this book.

| Reinforcement Learning | Game Theory |
| --- | --- |
| Agent | Player |
| Environment | Game |
| Reward | Utility |
| State | InfoSet / InfoState / State |
| Policy | Strategy |

Table 1.1: Reinforcement learning and game theory use different words for analogous concepts.

## 1.8 Structure

The thesis is structured into two parts. First — perfect information games — looks at search in perfect information games from a slightly unusual perspective. Looking at search from this perspective creates basic building blocks, concepts and ideas to be expanded to imperfect information games.

The second part — imperfect information games — is the core of the contributions of this thesis. Without the first part, the ideas and concepts would be accessible only to a very small community.

# Part I

# Perfect Information

# 2. Introduction

The first part of this thesis deals with perfect information games. The point of this section is not to introduce any new results or insights into search in perfect information games. Rather, it is to look closely at the individual search components — locality, decomposition and value functions. We will look at search from a particular perspective. And while this perspective is arguably not very useful for perfect information games, it will allow us to introduce search in imperfect information in a much more comprehensible way. It will also allow us to appreciate the different complications in the fundamental search components, as we will be extending these concepts in the second section of this thesis.

## 2.1   Chapters Structure

**Formal Models**   Chapter 3 introduces a simple game-tree formalism that is sufficient for the games in question.

**Offline Policies**   Chapter 4 motivates optimal policies. Introduces concepts of best response, maximin and Nash equilibria and discusses different ways a policy can be evaluated. It also shows the equivalency of maximin and Nash for two player zero-sum games. We finish with some dis-concerning notes on multiplayer settings.

**Sub-games**   Chapter 5 presents a key decomposition concept of sub-games.

**Offline Solve**   Chapter 6 discusses offline algorithms for producing optimal policies. We also introduce self-play style methods, including the popular approach o self-play via independent reinforcement learning.

**Online Settings**   Chapter 7 shows that algorithms operating in online settings need to be thought of and evaluated differently to offline algorithms. We show some counter-intuitive examples of an online algorithm's worst case performance. We then introduce an appropriate formalism and metric for the online setting, including a framework connecting the performance of online and offline algorithms.

**Search**   Chapter 8 discusses basic building blocks of search and how they come together. It introduces online minmax algorithm and includes well-known limited depth minmax search algorithm.

**Examples**   Chapter 9 lists some well known examples of successful utilization of search in perfect information games.

# 3. Formalisms

For perfect information games, we will settle with a single formalism for sequential decision making — theperfect information game tree.

## 3.1  Perfect Information Game Tree

A model that fits many popular perfect information games (e.g. chess, go, checkers) is a simple game tree (Figure 3.1). In this tree, the nodes correspond to game states (e.g. chess position) while the edges correspond to the legal actions (e.g. moving a queen).

**Definition 1.** *A perfect information game tree consists of:*

- *Set of players $\mathcal{N} = \{0, 1, 2\}$, where the player 0 is the chance player.*

- *Set of states (tree nodes) $\mathcal{S}$.*

- *Set of terminal states $\mathcal{Z} \subseteq \mathcal{S}$.*

- *Player to act in a state $s \in (\mathcal{S} \setminus \mathcal{Z})$ denoted as $p(s) \in \mathcal{N}$. As a single player acts in a state, we also define $\mathcal{S}_i = \{s \in S \,|\, p(s) = i\}$.*

- *The set of available actions (tree edges) for a state $s$, denoted as $\mathcal{A}(s)$.*

- *Mapping of state-action pair to a next state (child) $c(s, a) : \mathcal{S} \times \mathcal{A} \to \mathcal{S}$.*

- *Fixed strategy of a chance player $\pi_c : s \in \mathcal{S}_0 \to \Delta(\mathcal{A}(s))$.*

- *Return for a player i a terminal state $R_i : \mathcal{Z} \to \mathcal{R}$.*



Figure 3.1: Simple game tree.

## 3.1.1  Derived Concepts

We denote the child state after taking the action $a \in \mathcal{A}(s)$ in a parent $s$ as $sa$, and analogically for longer sequences $sa_1a_2 \ldots a_n$. We also use $a_i \sqsubseteq a_j$ to denote that state $a_i$ is a predecessor of $a_j$. Furthermore, $\Delta_R$ is the delta of maximim/minimum utilities in the game: $\Delta = \max_{z \in \mathcal{Z}} R_i(z) - \min_{z \in \mathcal{Z}} R_i(z)$

### 3.1.2  Policy

A behavioral policy is a mapping from a state to a distribution over the available actions $s \to \Delta(\mathcal{A}(s))$. We denote a policy of a player $i$ as $\pi_i$ and the policy in a state $s \in \mathcal{S}$ as $\pi_i(s)$. The set of all policies of player $i$ is then $\Pi_i$. As a single player acts in a state, we often use simply $\pi(s)$. Furthermore, the probability mass for an action $a \in \mathcal{A}(s)$ is $\pi_i(s, a)$. A policy profile consists of strategies of both player $\pi = (\pi_1, \pi_2)$. Given a policy profile, we also define a reach probability of a state $s \in \mathcal{S}$ as $P^\pi(s) = \prod_{s'a \sqsubseteq s} \pi(s', a)$. Finally, useful concept is the support (Definition 2)

**Definition 2** (Support). *Support of a policy $\pi_i$ is the set of actions with non-zero probability under that policy $supp^{\pi_i}(s) = \{a \in \mathcal{A} \,|\, \pi(s, a) > 0\}$.*

### 3.1.3  Expected Reward

Unlike in single agent settings, we need to know both players' policies to compute the reward. And since both the policies and the environment (chance player) can be stochastic, we care about the expected reward. As the reward is defined over the terminal states $\mathcal{Z}$, the expected reward of player $i$ given the strategy profile is simply (3.1).

$$R_i(\pi) = R_i(\pi_1, \pi_2) = \sum_{z \in \mathcal{Z}} P^\pi(z) R_i(z) \tag{3.1}$$

### 3.1.4  Averaging Policies

One must exercise caution during any averaging in the space of behavioral strategies. A common mistake is to simply average the strategy locally, in each individual state. This is incorrect as for $\pi_1^c = 0.5\pi_1^a + 0.5\pi_1^b$, we expect $R_1(\pi_1^c, \pi_2) = 0.5R_1(\pi_1^a, \pi_2) + 0.5R_1(\pi_1^b, \pi_2)$. Figure 3.2 shows a counterexample to this expectation for naive per-state averaging of the policy.

The issues is because the dynamics are inherently sequential — the distribution over the states $P^\pi(z)$ depends on the full sequence. Proper averaging in behavioral state must thus take into account the reach probability of the state (Figure 3.2d). Section 16.3.1 will introduce sequence form representation that allows for convenient linear operations directly on this strategy representation.

### 3.1.5  Constant-Sum and Zero-Sum Games

For a constant-sum game, the reward is strictly competitive — gain of one agent is the loss of the other. The rewards of both players in each terminal sums to a constant (Definition 3). This is a property of all win/lose/draw games (e.g. chess) as well as most games where money exchange hands (e.g. poker).

**Definition 3** (Constant Sum Game).

$$\exists c \in \mathcal{R} \,:\, \forall z \in \mathcal{Z} \,:\, \sum_{i \in \mathcal{N}} R_i(z) = c$$

(a) $\pi_1^a$      (b) $\pi_1^b$      (c) $\pi_1^c$ - naive average   (d) $\pi_1^c$ - proper average

Figure 3.2: Consider the highlighted state $s$ against an opponent who plays to reach that state, and its reach probability $P^\pi(s)$. a) $P^{\pi_1^a}(s) = 0.64$ b) $P^{\pi_1^b}(s) = 0.04$ c) Naive per-state averaging of the strategy simply averages the state strategy in isolation: $P^{\pi_1^c}(s) = 0.25 \neq 0.5 \cdot 0.64 + 0.5 \cdot 0.04$ d) Proper averaging of the strategy takes reach account into consideration:

A zero-sum game is then simply a game for which $c = 0$. There is practically no difference between constant-sum and zero-sum games, as the reward can simply be shifted to zero with no strategic or algorithmic implications. It is thus common to use the term zero-sum game even if the presented results hold true for any constant.

# 4. Optimal Policies

Formal definitions for the games and policies in hand, we now come back to the question of optimal policies. This time, we will be able to formally define these concepts and important properties. Namely, we start with a simple concept of best response that maximizes utility given a fixed opponent. Next we introduce a maximin — policies that optimize against the worst case scenario. The final concept is Nash equilibrium — a stationary strategy profile where no player has incentive to deviate. We then discuss some connections between these concepts. Importantly, we show that maximin policies and Nash equilibrium collapse for two player zero sum games. We finish with some notes and observations regarding multiplayer settings, where the concept of maximin and Nash equilibrium differ.

## 4.1 Best Response

Best response is a concept that will be used throughout this book. It is simply a policy that maximizes return against a fixed policy.

**Definition 4** (Best Response). *Best response against a policy $\pi_i$ is:*

$$\arg\max_{\pi_{-i} \in \Pi_{-i}} R_{-i}(\pi_i, \pi_{-i})$$

We use $\mathbb{BR}(\pi_i)$ to denote the set of best response policies against the policy $\pi_i$. Note that for zero-sum games, opponent maximizing their reward is equivalent to opponent minimizing our reward.

$$\arg\max_{\pi_{-i}} R_{-i}(\pi_i, \pi_{-i}) = \arg\min_{\pi_{-i}} R_i(\pi_i, \pi_{-i})$$

As this means the player's value against any best-response strategy is unique, we denote this unique value as $BRV_i(\pi_i)$.

$$BRV_i(\pi_i) = \min_{\pi_{-i}} R_i(\pi_i, \pi_{-i}) = -\max_{\pi_{-i}} R_{-i}(\pi_i, \pi_{-i})$$

### 4.1.1 Properties

How hard is it to compute a best response? The first look at the definition is not necessarily helpful from the computation perspective — we are taking an arg min over all possible policies. Fortunately, computing a best response strategy can be done via a single bottom-up pass traversal of the game tree (Algorithm 1). As the algorithm produces deterministic policies, we immediately get the Lemma 1

**Lemma 1.** *There is always a deterministic best response.*

## 4.2 Maximin

The worst-case reasoning (Section 1.3.1) motivates a policy that maximizes the reward against the worst case opponent. We are now able to formally introduce the maximin policy (Definition 5).

---

**Algorithm 1** Best Response

---

1: **function** BestResponseDFS($s \in S, i \in N$)
2:                             ▷ Terminal state, return terminal utility.
3:      **if** $s \in \mathcal{Z}$ **then**
4:          **return** $R_i(s)$
5:                    ▷ Compute the best response values of the children states.
6:      **for** $a : \mathcal{A}(s)$ **do**
7:          $v[sa] = \text{BestResponseDFS}(sa, i)$
8:
9:      **if** $\text{p}(s) = \text{i}$ **then**               ▷ Best response decision, taking max.
10:         $\pi(s) = \arg\max_{a \in \mathcal{A}(a)} v[sa]$
11:         **return** $\max_{a \in \mathcal{A}(a)} v[sa]$
12:     **else**                   ▷ Other player's decision, use their strategy
13:         **return** $\sum_{a \in \mathcal{A}(a)} \pi(s, a) v[sa]$
14:
15: **function** BestResponse ▷ Start the DFS computation in the initial state
16:     $\text{BestResponseDFS}(s_0, i)$

---

**Definition 5** (Maximin Policy). *Maximin policy of a player i is:*

$$\arg\max_{\pi_i \in \Pi_i} \min_{\pi_{-i} \in \Pi_{-i}} R_i(\pi_i, \pi_{-i}) = \arg\max_{\pi_i \in \Pi_i} BRV_i(\pi_i) \tag{4.1}$$

We denote the set of all such maximin policies for a player $i$ as $\mathbb{MAXIMIN}_i$.

$$\mathbb{MAXIMIN}_i = \{\pi_i \mid BRV_i(\pi_i) = \max_{\pi_i'} BRV_i(\pi_i')\} \tag{4.2}$$

Furthermore, the corresponding value of the maximin policy is the maximin value, denoted as $\underline{v_i}$ (Definition 6)

**Definition 6** (Maximin Value). *Maximin value is the value of the maximin policy:*

$$\underline{v_i} = \max_{\pi_i \in \Pi_i} \min_{\pi_{-i} \in \Pi_{-i}} R_i(\pi_i, \pi_{-i}) = \max_{\pi_i \in \Pi_i} BRV_i(\pi_i) \tag{4.3}$$

## 4.2.1   Minimax Theorem

We will now investigate the important relation between the players' respective maximin values $\underline{v_i}$ and $\underline{v_{-i}}$.

$$\underline{v_i} = \max_{\pi_i} \min_{\pi_{-i}} R_i(\pi_i, \pi_{-i}) \tag{4.4}$$

$$\underline{v_{-i}} = \max_{\pi_{-i}} \min_{\pi_i} R_{-i}(\pi_i, \pi_i) \tag{4.5}$$

First, we will re-phrase equation (4.5) in terms of $R_i$:

$$\underline{v}_{-i} = \max_{\pi_{-i}} \min_{\pi_i} R_{-i}(\pi_i, \pi_i) \tag{4.6}$$

$$= \max_{\pi_{-i}} \min_{\pi_i} -R_i(\pi_i, \pi_i) \tag{4.7}$$

$$= \max_{\pi_{-i}} - \max_{\pi_i} R_i(\pi_i, \pi_i) \tag{4.8}$$

$$= - \min_{\pi_{-i}} \max_{\pi_i} R_i(\pi_i, \pi_i) \tag{4.9}$$

Theorem 2 then states a critical result — the maximin values are in balance $\underline{v}_i = -\underline{v}_{-i}$. We refer to this unique value as the game value and denote it as $GV_i$.

**Theorem 2** (Minimax Theorem)**.**

$$\max_{\pi_i} \min_{\pi_{-i}} R_i(\pi_i, \pi_{-i}) = \min_{\pi_{-i}} \max_{\pi_i} R_i(\pi_i, \pi_i) \tag{4.10}$$

The minimax theorem, proven by John Von Neumann in 1928 (Neumann, 1928) has dramatic consequences for two player zero sum games (the proof is for both perfect and imperfect information games). Von Neumann himself later wrote "As far as I can see, there could be no theory of games on these bases without that theorem" (Von Neumann and Fréchet, 1953). We will prove the minimax theorem for the more general case of imperfect information in Section 16.2.3.

## 4.3   Nash equilibrium

Another key solution concept — Nash equilibrium — is based on a stationary property. The idea is that if for a strategy profile $(\pi_i, \pi_{-i})$, none of the players benefit by deviating from their policy, the profile forms an equilibrium (Definition 7).

**Definition 7** (Nash Equilibrium)**.** *Strategy profile $(\pi_i, \pi_{-i})$ forms a Nash equilibrium if none of the players benefit by deviating from their policy.*

$$\forall i \in N, \forall \pi_i' \; : \; R_i(\pi_i, \pi_{-i}) \geq R_i(\pi_i', \pi_{-i})$$

We denote the set of all such strategy profiles as $\mathbb{NEQ}$ — strategy profile $(\pi_1, \pi_2)$ forms a Nash equilibrium iff $(\pi_1, \pi_2) \in \mathbb{NEQ}$. Note that the $\mathbb{NEQ}$ can also be formulated using the best response (Lemma 3)

**Lemma 3.** *Strategy profile $(\pi_i, \pi_{-i})$ forms a Nash equilibrium iff*

$$\forall i \in N \; : \; \pi_i \in \mathbb{BR}(\pi_{-i})$$

*Proof.* Follows directly from the definition of $\mathbb{NEQ}$ and $\mathbb{BR}$. $\qquad\square$

## 4.4 Nash Equilibrium vs Maximin

It might be unclear why we introduced both maximin and Nash equilibria as solution concepts. While the properties discussed do sound intriguing, which solution concept should we prefer: the stationary property of Nash Equilibria, or the worst case reasoning of maximin? While the maximin strategy is a single strategy that optimizes against the worst case, Nash equilibria was introduced as a strategy profile (that is, a pair of strategies).

It turns out the these solution concepts in two player zero sum games are the same!

**Theorem 4** (Nash is maximin). *For two player zero-sum games:*

$$(\pi_1, \pi_2) \in \mathbb{NEQ} \implies \pi_1 \in \mathbb{MAXIMIN}_1 \wedge \pi_2 \in \mathbb{MAXIMIN}_2$$

*Proof.* By contradiction. Denote the worst case value of $\pi_i$ as $a = BRV_i(\pi_i)$. Note that since the $\pi_{-i}$ is a best response (Lemma 3) and $R_i = -R_{-i}$ (zero-sum game), it must be the case that $a = R_i(\pi_i, \pi_{-i})$. Let there be another policy $\pi_i^*$ that does strictly better in worst case, that is, $b = BRV_i(\pi_i^*)$ and $b > a$. Then the player $i$ could improve by switching to policy $\pi_i^* : R_1(\pi_i^*, \pi_{-i}) > R_1(\pi_i, \pi_{-i})$, contradicting the definition of Nash equilibrium. $\qquad\square$

**Theorem 5** (Minimax is Nash). *For two player, zero-sum games:*

$$\pi_1^* \in \mathbb{MAXIMIN}_1 \wedge \pi_2^* \in \mathbb{MAXIMIN}_2 \implies (\pi_1^*, \pi_2^*) \in \mathbb{NEQ}$$

*Proof.* Let $\underline{v_i}, \underline{v_{-i}}$ be the max min values of the respective players. We will use the fact that $\underline{v_i} = -\underline{v_{-i}}$ (Theorem 2) to show a contradiction. Assume that the strategy profile is not Nash and there exists $\pi_1'$ such that $u_1(\pi_1', \pi_2^*) > u_1(\pi_1^*, \pi_2^*)$:

$$u_1(\pi_1', \pi_2^*) > u_1(\pi_1^*, \pi_2^*) \geq \underline{v_i}$$
$$-u_2(\pi_1', \pi_2^*) = u_1(\pi_1', \pi_2^*) > u_1(\pi_1^*, \pi_2^*) \geq \underline{v_i}$$
$$\underline{v_i} = -\underline{v_{-i}} \geq -u_2(\pi_1', \pi_2^*) = u_1(\pi_1', \pi_2^*) > u_1(\pi_1^*, \pi_2^*) \geq \underline{v_i}$$

Which concludes the contradiction. $\qquad\square$

Putting Theorem 4 and Theorem 5 together results in Corollary 5.1.

**Corollary 5.1** (Maximin and Nash Collapse). *For two player zero sum games, solution concepts of maxmin and Nash equilibria are the same.*

$$\pi_1 \in \mathbb{MAXMIN}_1 \wedge \pi_2 \in \mathbb{MAXMIN}_2 \iff (\pi_1, \pi_2) \in \mathbb{NEQ} \qquad (4.11)$$

## 4.5 Properties of Optimal Policies

We now have formal understanding of the maximin and Nash equilibria solutions concepts. We also know that in our settings, these concepts are the same. But what motivation do we have to follow such policies?

### 4.5.1 Guaranteed Value

Following an optimal policy guarantees the maximin game value against any opponent. Furthermore, Theorem 2 tells us that the maximin value of a player is in balance with the maximin value of the opponent $\max_{\pi_i} \min_{\pi_{-i}} R_i(\pi_i, \pi_{-i}) = -\max_{\pi_i} \min_{\pi_{-i}} R_{-i}(\pi_i, \pi_{-i})$. We thus directly get the important Corollary 5.2. This appealing property of optimal policies is of a great importance, and easily motivates such policies.

**Corollary 5.2.** *If we follow an optimal policy when playing as either player, the expected utility against any opponent is greater or equal to zero:*

$$(\pi_i, \pi_{-i}) \in \mathbb{NEQ} \,:\, R_i(\pi_i, \pi'_{-i}) + R_{-i}(\pi'_i, \pi_{-i}) \geq 0 \,\forall \pi'_i, \pi'_{-i}$$

Corollary 5.2 essentially tells us that if we follow an optimal policy, we can't lose. If we face an optimal opponent, we draw. If we face a sub-optimal opponent, we might win (i.e. expect a positive reward), although we are certainly not guaranteed to. This is because even if the opponent follows a highly exploitable strategy, the optimal policy we follow might not take advantage of any of the mistakes. Even though this can easily happen in perfect information games, we defer examples and more discussion of this behavior to imperfect information games (Section 14.5).

### 4.5.2 Convexity

The optimal policy guarantees us the game value regardless of the opponent we face. And while there might be multiple optimal strategies, it does not matter which one we follow — we still enjoy the same guarantees. But how many optimal policies are there? What structure does the set of all optimal policies have? It is easy to verify that any convex combination[1] of optimal policies results in an optimal policy.

**Lemma 6.** *The set of all Nash equilibria $\mathbb{NEQ}$ is convex.*

*Proof.* We defer this result to the more general case of imperfect information games, proven in Section 16.2, Theorem 20. □

**Corollary 6.1.** *There is either one unique Nash equilibrium, or infinitely many.*

*Proof.* Follows directly from Lemma 6. □

---

[1]One must be careful not to naively combine the behavioral strategies, proper combination as described in Section 3.1.4 is needed

Figure 4.1: In multiplayer games, Nash equilibria and maximin solution concepts differ. Action selected by a policy are denoted by the bold red line. a) Nash strategy profile, where the individual strategies are not maximin. b) Pair of maximin strategies that do not for Nash strategy profile.

## 4.6 Multiplayer Games

While the properties of Nash equilibria do sound very intriguing, they are lost once we leave the comfort of two player zero-sum games. First, observe that introducing more players is essentially the same as removing the zero-sum property. For non-zero sum games, one can always add a dummy player with no actions with the surplus reward.

### 4.6.1 Nash is not Maximin

Unlike in two player zero sum games, the solution concepts of Nash and maximin are different. See Figure 4.1 for simple counterexamples.

### 4.6.2 Existence

Under what circumstances do optimal policies exist? It is easy to believe that maximin strategy always exists as it is simply maximizing against a best response value. But is it also the case for Nash? It turns our that the existence is guaranteed for a very wide class of games — even for games with multiple players. Essentially, one can show that Nash equilibrium is guaranteed to exist for any finite game (Nash et al., 1950; Nash, 1951).

### 4.6.3 Unique Value

There is no longer a unique value of the game, i.e. the minimax theorem no longer holds. While in two player settings, any Nash equilibrium guarantees the same game value and thus the players could arbitrarily select any Nash to follow, this does not hold in multiplayer settings. Even worse, it could very well be that $(\pi_1^a, \pi_2^a) \in \mathbb{NEQ}, (\pi_1^b, \pi_2^b) \in \mathbb{NEQ}$ and both strategy profiles produce high reward for the players, but if they do not coordinate on which equilibrium to follow and they select different ones — they end up with $(\pi_1^a, \pi_2^b)$ where the value can be arbitrarily bad for both. This issue is also referred to as the equilibrium selection problem, and we defer an example of this behavior to Section 14.6.

### 4.6.4 Complexity

One can easily argue that maximin is as hard in multiplayer settings as it is for two players. We can just let a single opponent control all the other players.

For Nash equilibria, the question of complexity in multiplayer settings is more complicated. We can not use the same argument as for maximin, since Nash equilibria requires that none of the players benefits by deviating from the strategy profile — that is, each player individually. But the players could benefit from deviating if multiple agents deviate at once. In general, the problem is known to be PPAD-complete (Papadimitriou, 1994; Chen and Deng, 2006; Daskalakis et al., 2009).

### 4.6.5 What to Do

But maybe it is just the case that there is another, simple and intriguing solution concept for multiplayer games? Unfortunately, it seems that in general, things are quite gloomy. There seems to be no consensus of what solution concept would that be, and the results are mostly negative. The most promising line of work now focuses on online settings and coordinated strategies rather than on following a fixed strategy.

## 4.7 Evaluation

Given a sub-optimal policy $\pi_i$, we often want to evaluate how "close" to an optimal policy it is. As the worst-case value is no more than the game value: $\delta(\pi_i) = GV_i - BRV_i(\pi_i)$, common measure then is $NASHCONV(\pi) = \sum_i \delta_i$ and exploitability $= \frac{NASHCONV}{|\mathcal{N}|}$. Further popular concept is $\epsilon$-Nash equilibrium, where the players receive at most $\epsilon$ by deviating from the strategy profile: $\max_i \delta_i(\pi_i) \leq \epsilon$. Exploitability and $\epsilon$-optimal policies allow to objectively and quantitatively evaluate a policy, and contrasts its worst-case performance to the worst-case performance of an optimal policy.

But exploitability itself does not tell the full story about the strength of an agent. Strong chess agent that is beatable by a particularly clever line of play is a arguably a better chess player than an agent that always resigns. Yet, both of these agents have the same exploitability.

Sometimes, we are not interested in the worst-case performance, but rather in the qualitative performance of the agent in head to head evaluation against a specific pool of the opponents. The issue is that such performance strongly depends the players in the pool, but there are also inherent intransitivies in head-to-head evaluation. Given three players $\pi^a, \pi^b, \pi^c$, it could very well be that $\pi^a$ beats $\pi^b$, $\pi^b$ beats $\pi^c$ and $\pi^c$ beats $\pi^a$.

### 4.7.1 Approximate Evaluation

To compute exploitability or nashconv, we need to compute the exact best response. This requires a full tree traversal (Algorithm 1), which can quickly become intractable in large games. By fixing policy of one of the players, the game

collapses into a single agent environment (Bowling, 2003), where the optimal policy corresponds to a best response. One can thus opt for any standard reinforcement learning methods to learn (or approximate) this optimal best-responding policy (Greenwald et al., 2017) and use the resulting approximate policy to define a corresponding metric — approximate best response (Timbers et al., 2020). Local best response (Section 18.1.2) is then another method approximating the best response, combining search and heuristic value functions (Lisy and Bowling, 2017).

# 5. Sub-games



Figure 5.1: a) Current state to reason about. b) All the states relevant for the current decision form a sub-game.

As the name suggests, a sub-game is a sub-problem of a game. Sub-games are a fundamental concept for online search methods, and there are two critical ways the are used. While both cases essentially decompose the problem and combine the solutions, they are very much semantically distinct.

First, sub-games allow us to reason about a small sub-problem in isolation. Search methods reason about sub-games, namely about the sub-games defined by the currently observed state.

Second, sub-games are the basis of value functions. If the sub-game defined by the current state is still too large, it allows for further decomposition. Rather than reasoning about all future states in the sub-game, we can look forward only some number of steps and use the value of the future states (future sub-games).

## 5.1  Perfect Information Sub-game

Given a current state, what is the smallest set of states relevant to the current decision? In perfect information games, it is simply the state and all the future states, a sub-tree rooted in that state (Figure 5.1). To verify that is the case, note that a policy for these states is sufficient to compute a value for the current state, as the policies in all other states are irrelevant.

The defining property of sub-game is that — as the name suggests — it is a game. It needs to be well-defined game to be reasoned about in isolation. Fortunately, this is again straightforward in perfect information games formalism. All we have to do is to take a sub-tree rooted in a state (Definition 8).

**Definition 8.** *Sub-game of a states $s \in \mathcal{S}$ is a game formed by the sub-tree of $s$.*

## 5.2 Sub-game Values

As sub-game forms a well-defined game, we can readily define the value of a sub-game (Definition 9).

**Definition 9** (Sub-game Value). *Sub-game value is the game value of the corresponding game.*

Sub-game values will turn out to be very important for decomposition and search methods, as they will allow allow us to reason about optimal policies in a state by reasoning about future states and the respective sub-game values.

## 5.3 Sub-game Perfect Optimal Strategy



Figure 5.2: a) Optimal policy that guarantees a game value of 1, but is not sub-game perfect as it does not play optimally in the highlighted sub-game. b) Sub-game perfect optimal policy.

With the notion of sub-games in hand, we can define a useful refinement of the optimal policies. A strategy in the full game can be optimal and guarantee the game value while playing sub-optimally in some of the sub-games. An important refinement of the optimality concept is then sub-game perfect strategy, which is optimal in each sub-game (Definition 10). See Figure 5.2 for an example.

**Definition 10** (Sub-game Perfect Strategy). *A sub-game perfect strategy is optimal for each sub-game $s \in \mathcal{S}$.*

# 6. Offline Solving

As previous chapters motivated the policies we are interested in, we will now have a look at methods that allow us to find the desired strategies. The setting we start with is that of offline solving. That is, the algorithm computes the strategy for all the possible decision points prior to game play. While we will be mostly concerned with tabular algorithms and representations, we also discuss some approximation methods.

## 6.1   Minimax Algorithm



Figure 6.1: Minimax algorithm traverses the game tree bottom-up, producing optimal policies in the sub-games and sending the sub-game values up the tree.

The minimax algorithm is a very simple method of computing an optimal policy for a game tree. This algorithm is often introduced as producing "optimal policies" without any deeper elaboration on what that really means. Previous chapters allow us to get more insights into this algorithm, as we now formally understand optimal policies and their properties (Chapter 4) as well as sub-games (Chapter 5).

The minimax algorithm traverses the game tree bottom-up, producing optimal policies in the sub-games and sending the sub-game values up the tree. These sub-game values are then turned into optimal policies for larger sub-games until the process reaches the root state (Figure 6.1 and Algorithm 2). The minimax algorithm is a great example of decomposition, as it inherently relies on the notion of sub-games and sub-game values.

Note that as we are computing the policies for both players in the bottom-up fashion, we are not using any knowledge of how the players play further up the tree. This is because in perfect information, sub-games require no such information as the sub-game is defined simply as a sub-tree. This is yet another property that does not hold in imperfect information as we will see in Chapter 15.

### 6.1.1   Resulting Policy

To verify that Algorithm 2 produces the desired solution, notice that the resulting policies best-respond to each other (see the best response Algorithm 1). Lemma 3 then guarantees that such policy pair forms a Nash equilibrium. Finally thanks to Corollary 5.1, Nash equilibrium and maximin collapse in our settings.

### 6.1.2   AlphaBeta Pruning

While the minimax algorithm traverses the entire tree, there is a way to potentially traverse only part of the tree while keeping the same guarantees. Alpha-beta

pruning allows for potentially visiting only $\sqrt{|\mathcal{S}|}$ of nodes rather than $|\mathcal{S}|$ (Knuth and Moore, 1975).

---

**Algorithm 2** Minimax

---

1:
2: **function** MINIMAX($s \in S$)
3:      **if** $s \in Z$ **then return** $R_1(s)$                         ▷ Terminal state.
4:
5:      **for** $a \in \mathcal{A}$ **do**
6:          $v[sa] = $ MINIMAX($sa$)
7:
8:      **if** p(s) $= 1$ **then**                      ▷ Maximizing player to act.
9:          $\pi(s) = \arg\max_{a \in \mathcal{A}(s)} v[sa]$
10:          **return** $\max_{a \in \mathcal{A}(s)} v[sa]$
11:      **else if** p(s) $= 2$ **then**                  ▷ Minimizing player to act.
12:          $\pi(s) = \arg\min_{a \in \mathcal{A}} v[sa]$
13:          **return** $\min_{a \in \mathcal{A}(s)} v[sa]$
14:      **else**                                 ▷ Chance player to act.
15:          **return** $\sum_{a \in \mathcal{A}(s)} \pi(s, a) v[sa]$

---

## 6.2 Self-Play Style Methods

Self-play is a powerful idea where agents improve their policies in a self reinforcing loop (Figure 6.2). The hope is that as the strategies keep improving, they will eventually converge to the optimal ones. In the self-play framework, agents produce a sequence of strategies $(\pi_1^t, \pi_2^t)$ and algorithms differ in important details. First is how the algorithm produces the next strategy in the sequence. Second is which strategy is expected to converge (e.g. the average one or the current one).

### 6.2.1 Self Play via Independent Reinforcement Learning

Self-play via independent reinforcement learning is a popular method coming from the reinforcement learning community. This approach simply uses a single-agent reinforcement learning method to improve the agent's policy against the opponent. The naive approach is to simply run a single agent reinforcement learning method of choice independently for both players, each improving against the last strategy of the opponent. This is indeed at the very heart of many well-known success stories of reinforcement learning methods in two player games (Tesauro, 1995; Silver et al., 2017b). The reasoning sounds appealing — if we keep improving, surely we will end up with an optimal policy?

**Convergence**

What are the convergence guarantees of this method? To simplify the analysis, we can let the reinforcement learning method fully converge during each self play episode. We know from the Section 4.7.1 that this produces the best response.

Figure 6.2: Self-play methods produce a sequence of strategies that improve against each other. As fixing one agent makes the environment single-agent, the environment becomes MDP during each improvement step.



Figure 6.3: Self play methods that simply best respond to the last policy of the opponent do not necessarily converge.

We thus end up in a self-play sequence where at each time step, the policy best-responds to the previous policy of the opponent: $\pi_i^t \in \mathbb{BR}_i(\pi_{-i}^{t-1})$. We refer to this sequence as the best responding sequence.

Lemma 7 shows that if such sequence converges, the resulting strategy profile is optimal.

**Lemma 7.** *If the best responding sequence converges* $\exists t \forall t' > t : \pi^{t'} = \pi^t$, *the strategy profile* $\pi^t$ *forms a Nash equilibrium.*

*Proof.* Since the sequence is best-responding $\pi_i^{t'} \in \mathbb{BR}(\pi_{-i}^{t'-1})$ and due to convergence $\pi_i^{t'} = \pi_i^{t'-1}$. Thus the policies $\pi_i^{t'}$ and $\pi_{-i}^{t'}$ best respond to each other, forming a Nash equilibrium (Lemma 3). $\qquad\square$

But does this method necessarily converge? While it is tempting to say yes, the answer is no — even in perfect information games! Consider a game in Figure 6.3 and the following sequence:

$$(\pi_1^t, \pi_2^t) = \begin{cases} (L, ad), & \text{if } t \% 2 = 0 \\ (R, bc), & \text{otherwise} \end{cases}$$

.

The sequence satisfies the best-responding property, but it never converges. This would not be bad if the sequence would simply alternate between different optimal policies. But at each time step, $\pi_2^t$ is highly exploitable. This sequence does not converge, and produces a highly exploitable current strategy as well as the average strategy.        Further analysis of the sequence reveals that a necessary ingredient in this counter-example is that the best response of the second player is never sub-game perfect (Section 5.3). A best responding sequence is indeed guaranteed to produce optimal policies if the best response is sub-game perfect. Note that this is true only in perfect information games and we will see in Section 16.4 that these approaches fail in imperfect information.

# 7. Online Settings

Hopefully by this point, the reader is convinced of the importance of Nash equilibria. We have learned about its properties and introduced methods producing this optimal policy. So far, we dealt with offline policies — prior to playing the game, we would compute a policy for each state and then store it. During the actual game-play, we would simply follow this policy.

But search algorithms operate in fundamentally different settings, producing a strategy for a state only once it is visited. This online setting requires a different and careful analysis, as the offline concepts do not apply. Unlike offline policies, online algorithms can condition the computation on the past experience and games, allowing for opponent adaptation (Bard, 2016) or re-using parts of the previous computation (Silver et al., 2016). In general, an online algorithm can produce game dynamics that is not consistent with any offline algorithm (Lemma 8).

We will introduce analogous concepts to $\epsilon$-Nash equilibria, $\epsilon$-Soundness, generalizing the concept of Nash equilibria to online settings. Sound algorithms guarantee (in expectation) the game value against arbitrary opponent. While $\epsilon$-soundness is the proper measure of worst-case performance, it can be hard to compute for some online algorithms as we might need to construct an exponentially large response game. We therefore introduce a consistency hierarchy that allows us to formalize when on online algorithm plays consistently with some Nash equilibrium. The strongest of such presented hierarchies, strong global consistency, then allows one to compute exploitability in a particularly easy way using tabularization.

The results presented in this Chapter (and the following results for imperfect information case in Chapter 19 ) are one of the contributions of this thesis.

## 7.1   Cautionary Example

To further illustrate the deceiving nature of online settings, consider the following algorithm. The algorithm internally uses an offline algorithm that produces provably optimal strategy. During play, the online algorithm i) runs the offline algorithm for some more time ii) retrieves a strategy $\pi_i(s)$ and takes an action $a \sim \pi_i(s)$. As the algorithm always "follows" an equilibrium, we would naturally expect this algorithm to be sound. This is not necessarily the case and such an algorithm can potentially be highly exploitable (Theorem 9). See Section 7.7.1 for a concrete example. Imperfect information games (Chapter 19 ) will then bring even more challenging examples.

## 7.2   Repeated Game

To properly analyze the performance of an algorithm, we need to introduce the repeated game. A repeated game consists of a sequence of individual matches (e.g. playing some number of chess games against the world champion). As a match progresses, the algorithm produces a strategy for a visited state on-line,

once it actually observes the state. We are then interested in the accumulated reward of the agent during the span of the repeated game. Of particular interest will be the expected reward against a worst-case adversary.

Formally, the repeated game $p$ consists of a finite sequence of $k$ individual matches $p = (z_1, z_2, \ldots, z_k)$, where each match $z_i \in \mathcal{Z}$ corresponds to the sequence of states and actions leading to that terminal state.

## 7.3 Online Algorithm

An online algorithm $\Omega$ maps a state visited during a repeated game to a strategy, while possibly using and updating its internal state (Def. 11). As the state of the algorithm is a function of previously visited (queried) states, we can also use $\Omega^{(z_1, \ldots, z_{k-1})}(s)$ to denote its output in a particular internal state corresponding to past experience.

**Definition 11.** *(Šustr et al., 2020, Definition 2) Online algorithm $\Omega$ is a function $s \times \Theta \mapsto \Delta(\mathcal{A}_i(s)) \times \Theta$ that maps an information state $s \in \mathcal{S}$ to a strategy $\Delta(\mathcal{A}_i(s))$, while possibly making use of algorithm's state $\theta \in \Theta$ and updating it. We denote the algorithm's initial state as $\theta_0$.*

### 7.3.1 Stateless vs Stateful

A special case of an online algorithm is a stateless algorithm, where the output of the function is independent of the algorithm's state. If the output depends on the algorithm's state, we say the algorithm is stateful algorithm. The distinction between stateful and stateless algorithms has some important consequences, as the state is the "memory" of the algorithm allowing any conditioning on the past.

### 7.3.2 Reward

Given two online players $\Omega_1, \Omega_2$, we use $P^k_{\Omega_1, \Omega_2}$ to denote the distribution over all the possible repeated games $p$ of length $k$ when these two players face each other. The average reward of $p$ is then $R_i(p) = \frac{1}{k} \sum_{j=1}^{k} u_i(z_j)$. Finally, we use $\mathbb{E}_{p \sim P^k_{\Omega_1, \Omega_2}}[R_i(p)]$ to denote the expected average reward when the players play $k$ matches.

While nothing stops the online algorithm from simply following a fixed policy profile, online settings allow for more general dynamics (Lemma 8).

**Lemma 8.** *An online algorithm $\Omega$ can produce game dynamics $P^k_\Omega$ that no offline strategy can.*

*Proof.* Consider a simple game with a single decision state $s_1$, $\mathcal{A}(s_1) = \{L, R\}$ and the following stateful online algorithm with state $\theta \in \mathbb{Z}$:

$$\Omega(s_1, \theta) = \begin{cases} (L, \theta + 1), & \text{if } \theta \% 2 = 0 \\ (R, \theta + 1), & \text{otherwise} \end{cases}$$

$\square$

## 7.4 Soundness

We are now ready to formalize the online concept analogous to Nash equilibrium. Exploitability / $\epsilon$-equilibrium considers the expected utility of a fixed strategy against a worst-case adversary in a single match. The analogous concept for the online settings in a repeated game is $\epsilon$-soundness (Definition 12). Intuitively, an online algorithm is $\epsilon$-sound if and only if it is guaranteed the same reward as if it followed a fixed $\epsilon$-equilibrium.

**Definition 12.** *(Šustr et al., 2020, Definition 4) For an $\epsilon$-sound online algorithm $\Omega$, the expected average reward against any opponent is at least as good as if it followed an $\epsilon$-Nash equilibrium fixed strategy $\pi$ for any number of matches $k$:*

$$\forall k \forall \Omega_2 \;:\; \mathbb{E}_{p \sim P^k_{\Omega,\Omega_2}}[R(p)] \geq \mathbb{E}_{p \sim P^k_{\pi,\Omega_2}}[R(p)]. \tag{7.1}$$

## 7.5 Response Game



Figure 7.1: To compute the $\epsilon$-soundness, we construct a $k$-repeated game where we replace the decisions of the online algorithm $\Omega$ with a fixed chance policy $\pi_0$, resulting in a single-agent game of exponential size. In the presented example, $k = 2$.

To compute the $\epsilon$-soundness as in Definition 12, we need construct a repeated game, where we replace the decisions of the online algorithm with stochastic (chance) transitions. As we allow the online algorithm to be stateful and thus produce strategies depending on the game trajectory, the response game must also reflect this possibility. The resulting game is thus exponential in size as it reflects all possible trajectories of $k$ matches. The chance policy $\pi_0$ for a state corresponding to past experience $p = (z_1, z_2, \dots)$ and current state $s$ is then $\pi_0(s^{p.s}) = \Omega^p(s)$. We call this single-player game a $k$-step response game (Figure 7.1).

## 7.6 Tabularized Offline Strategy

When an online algorithm is stateless and thus produces the same strategy for an information state regardless of the previous matches, there is no need for the

Figure 7.2: Locally consistent algorithm can be highly exploitable. Consider two Nash equilibria: $\pi_1 = \{s_1 : L, s_2 : Y\}, \pi_2 = \{s_1 : R, s_2 : X\}$. Online algorithm $\Omega = \{s_1 : L, s_2 : Y\}$ is locally consistent but highly exploitable.

$k$-response game. A fixed strategy $\pi$ sufficiently describes the behavior of the algorithm and the notions of exploitability and soundness collapse. To compute the corresponding fixed strategy, one simply queries the online algorithm for all the information states in the game $\pi(s) = \Omega(s) \forall s \in \mathcal{S}$. We refer to this strategy as the tabularized strategy.

## 7.7 Search Consistency

To prove that an online search algorithm is $\epsilon$-sound, we often want to formally state that the online algorithm plays "consistently" with an $\epsilon$-equilibrium. This allows one to directly bound the $\epsilon$-soundness of the online algorithm.

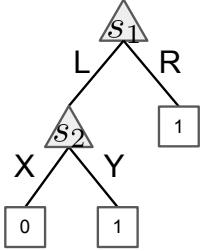We introduce three levels of consistency, with varying connections of how closely the online algorithm plays "just like" an $\epsilon$-equilibrium.

### 7.7.1 Local Consistency

The weakest of the connections, local consistency simply guarantees that every time we query the online algorithm, there is an $\epsilon$-equilibrium consistent with the produced behavioral strategy for that state (Definition 13).

**Definition 13.** *(Šustr et al., 2020, Definition 6) Algorithm $\Omega$ is locally consistent with $\epsilon$-equilibria if*

$$\forall p = (z_1, z_2, \ldots, z_k) \forall s \sqsubset z_k \exists \pi \in \mathbb{NEQ} : \Omega^{(z_1, \ldots, z_{k-1})}(s) = \pi(s). \qquad (7.2)$$

It might seem that local consistency is sufficient as it plays just like an equilibrium. Figure 7.2 presents a counterexample to this intuition. The example game includes two states $(s_1, s_2)$ for the first player and gives an example of two Nash equilibria $(\pi_1^1, \pi_1^2)$ that are in a way "incompatible". If an online algorithm follows $\pi_1^1$ and $\pi_1^2$ in $s_1$ and $s_2$ respectively, the resulting strategy can be highly exploitable (Theorem 9).

**Theorem 9.** *(Šustr et al., 2020, Theorem 7) An algorithm that is locally consistent might not be sound.*

Looking at the counterexample in more detail, one can notice that it relies on the fact that one of the Nash strategies is not sub-game perfect. Local consistency is indeed sufficient in perfect information games if the algorithm is consistent with a sub-game perfect equilibrium (Theorem 10). Theorem 10 does not hold in imperfect information settings and we will present a counter-example in Section 19.4.

**Theorem 10.** *(Šustr et al., 2020, Theorem 8) In perfect information games, an algorithm that is locally consistent with a subgame perfect equilibrium is sound.*
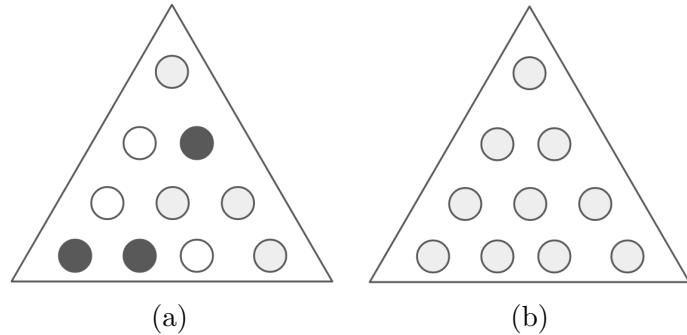
$$(a) \qquad\qquad\qquad (b)$$

Figure 7.3: (a) Local consistency — solving sub-games independently can lead to solution where policies in all individual states are consistent with an optimal policy, but each with possibly a different one. For this figure, the different colors represent different optimal policies for the full game. (b) Strong global consistency.

## 7.7.2 Global Consistency

Local consistency guarantees consistency for individual states in isolation. The problem we have seen is that the combination of behavioral strategies of individual states can yield highly exploitable strategy as there might be no equilibrium consistent with the resulting tabularized strategy. A natural extension is then to guarantee consistency for all the visited states in combination. Global consistency guarantees the existence of an equilibrium consistent with all the states visited during the gameplay (Definition 14). Unfortunately, even this stronger notion of consistency does not guarantee soundness (Theorem 11).

**Definition 14** (Global Consistency). *(Šustr et al., 2021, Definition 9) Algorithm $\Omega$ is globally consistent with $\epsilon$-equilibria if*

$$\forall p = (z_1, z_2, \ldots, z_k) \, \exists \pi \in \mathbb{NEQ} \, \forall s \sqsubset z_k \, : \, \Omega^{(z_1, \ldots, z_{k-1})}(s) = \pi(s). \qquad (7.3)$$

**Theorem 11.** *(Šustr et al., 2020, Theorem 10) An algorithm that is globally consistent with an $\epsilon$-equilibria might not be $\epsilon$-sound.*

But what if the algorithm keeps on playing the repeated game? While the global consistency with an equilibrium does not guarantee soundness, it guarantees that the average reward eventually converges to the game value in the limit (Theorem 12).

**Theorem 12.** *(Šustr et al., 2020, Theorem 11) For an algorithm $\Omega$ that is globally consistent with an $\epsilon$-equilibria:*

$$\forall k \; \forall \Omega_2 \, : \, \mathbb{E}_{p \sim P^k_{\Omega, \Omega_2}} [R(p)] \geq GV - \epsilon - \frac{|\mathcal{S}_1| \Delta_R}{k}, \qquad (7.4)$$

## 7.7.3 Strong Global Consistency

Essentially, the problem with global consistency is that it guarantees the existence of a consistent equilibrium *after* the game-play is generated. Strong global

Figure 7.4: (a) Possibly unknown offline policy $\pi$ the online search is re-solving. (b) Online search that is strong globally consistent with $\pi$.

consistency additionally guarantees that the game-play *itself* is generated consistently with an equilibrium. In other words, the online algorithm simply exactly follows a predefined equilibrium.

**Definition 15.** *(Šustr et al., 2020, Strong Global Consistency) Algorithm $\Omega$ is strongly globally consistent with $\epsilon$-equilibria if*

$$\exists \pi \in \mathbb{NEQ} \, \forall p = (z_1, z_2, \ldots, z_k) \, \forall s \sqsubset z_k \; : \; \Omega^{(z_1, \ldots, z_{k-1})}(s) = \pi(s). \qquad (7.5)$$

Strong global consistency guarantees that the algorithm can be tabularized, and the exploitability of the tabularized strategy matches $\epsilon$-soundness of the online algorithm.

## 7.8 Search as Re-Solving

The easiest way to argue that an online algorithm is sound is to make sure it is strongly globally consistent with some policy $\pi$. Such a search method then plays just like an offline policy $\pi$ would, except there is no explicit representation of $\pi$. Given a state $s$, we do not just solve for an optimal policy for that state. We rather re-solve a policy for that state, making sure it matches the "original" solve $\pi$. Search is then re-solving this policy step by step for all the visited states (Figure 7.4).

# 8. Search

We will now look into a particular variant of online algorithms — search. The algorithm is online as it computes the strategy for the current state $s \in \mathcal{S}$ only once the state is observed during game play. It then computes the strategy for the current state by searching forward, reasoning about the current sub-game.

We start with full lookahead settings, where the online algorithm constructs the full sub-game to be reasoned about. We describe the online minimax algorithm that computes a policy for the current state by running the minimax algorithm (Section 6.1) on the sub-game rooted in the current state. We then extend this algorithm to limited lookahead, where the search is truncated after some number of moves and the game value of the future sub-games is used as a substitute for continued search. Section 7.8 then discusses how to view these search algorithms in the context of re-solving, consistency framework an soundness.

We finish with Section 8.4 that includes several approaches for representing (approximate) value functions and Section 8.5 that includes discussion on the importance of size and structure of the lookahead tree.

## 8.1 Full Lookahead Online Minimax Algorithm

Full lookahead online minimax algorithm is a stateless online methods that produces a strategy for the current state by searching forward and reasoning about the sub-game rooted in the current state. This allows search to focus on the relevant part for the current decision, as sub-games are the smallest sub-problem we can reason about in isolation (Section 5). Full lookahead online minimax algorithm simply runs the minimax algorithm on the full sub-game rooted in the current decision $s \in \mathcal{S}$. It then produces a strategy for the current state $\Omega(s)$ and disregards the strategies of the other states of that sub-game.

The upside of this search algorithm is that if the current state is near the end of the game, we only need to reason about a small sub-game. The downside is that if the state is near the beginning of the game, we have to traverse the entire game tree anyway (Figure 8.1a). Thus the algorithm might be no better than the offline minimax algorithm, although it can avoid the memory needed to store the complete policy. While the algorithm is arguably not very practical, it is an important building block for the future algorithms and methods presented in this thesis.

## 8.2 Limited Lookahead Online Minimax Algorithm

The full lookahead online minimax algorithm considers all the states in the sub-game, looking forward until the end of the game. The idea of limited lookahead is to rather look forward only some number of steps, truncate the search and use the value (function) of the future sub-games (future sub-games of the current sub-game, Figure 8.1) in place of the continued search. The core of the idea is to

realize that the minimax algorithm only uses values of the future states to derive a strategy. Given the game values (value under an optimal policy) of sub-games, we are thus able to still derive an optimal policy for the current state.

The limited lookahead version of the online minimax algorithm builds a limited lookahead tree and runs the minmax algorithm within that tree. As we still need to know the value of the states at the end of the lookahead, the algorithm evaluates these states using a value function (Algorithm 3).



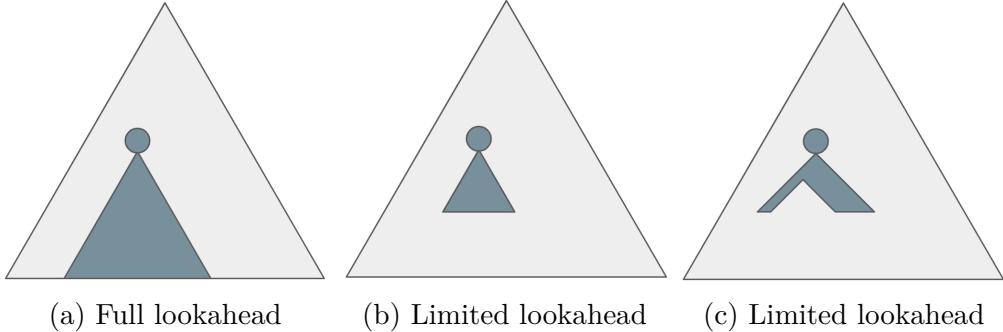(a) Full lookahead      (b) Limited lookahead      (c) Limited lookahead

Figure 8.1: (a) Full lookahead consists of the full subgame rooted in the current state. (b) and (c) Limited lookahead consists only of some of the future state. Size and shape of the lookahead can vary.

---

**Algorithm 3** Min Max Limited Lookahead

1:
2: **function** MinMaxLookaheadDFS($s \in \mathcal{S}$, $V : \mathcal{S} \to \mathcal{R}$)
3:      **if** $s \in Z$ **then return** $R_1(s)$              ▷ Terminal state.
4:      **if** $lookahead\_end(s)$ **then**
5:          **return** $V(s)$          ▷ Evaluate the state using value the function.
6:
7:      **for** $a \in \mathcal{A}$ **do**
8:          $v[sa] = $ MinMaxLookaheadDFS($sa$, $V$)
9:
10:      **if** p(s) = 1 **then**              ▷ Maximizing player to act.
11:          $\pi(s) = \arg\max_{a \in \mathcal{A}(s)} v[sa]$
12:          **return** $\max_{a \in \mathcal{A}(s)} v[sa]$
13:      **else if** p(s) = 2 **then**          ▷ Minimizing player to act.
14:          $\pi(s) = \arg\min_{a \in \mathcal{A}} v[sa]$
15:          **return** $\min_{a \in \mathcal{A}(s)} v[sa]$
16:      **else**                  ▷ Chance player to act.
17:          **return** $\sum_{a \in \mathcal{A}(s)} \pi(s, a) v[sa]$
18: **function** MinMaxLookahead($s \in \mathcal{S}$, $V : \mathcal{S} \to \mathcal{R}$)
19:      MinMaxLookaheadDFS($s$, $V$)
20:      **return** $a \sim \pi(s)$

---

## 8.2.1 Decomposition

There are two places where we just used the decomposition to sub-games. First,

the search methods produce a strategy for the current state by reasoning about the sub-game rooted in the current state. Second, sub-games are the basis of the value functions in the case of limited lookaheads. Conceptually similar decomposition will then take place in imperfect information.

## 8.3   Re-Solving and Consistency

Chapter 7 detailed the importance of non-locality. While search reasons inherently locally, it needs to make sure the resulting strategy is strongly globally consistent with an optimal strategy for the full game. Fortunately, this is particularly easy in this case, as the algorithm produces exactly the same behavioral strategies as the offline minmax algorithm. Another possible argument is to use Theorem 10.

Section 7.8 then discussed that in the case of strong global consistency, an online algorithm can be thought of as if it was re-solving an offline policy. In the case of online minimax algorithm, the policy being re-solved is the policy that the offline minimax algorithm (Algorithm 2) would produce if we were to run it on the full game (Figure 8.2).



Figure 8.2: a) Re-solving policy for the visited states. b) Online minimax full lookahead. c) Online minimax limited lookahead.

## 8.4   Approximate Value Functions

Exact value functions allow the search procedure to reason in a limited lookahead, as they provide the exact game value for the relevant future states at then end of the lookahead. But to compute the game value, we need to solve sub-games corresponding to these states anyway. Maybe instead of using the exact value, we can estimate the value of the future states.

Historically, the first approaches for approximating values were based on game specific heuristics (e.g. material in chess). Another popular method is to use sampling and estimate the value of a state by sampling some trajectories (Kocsis and Szepesvári, 2006; Browne et al., 2012).

Modern approaches learn the evaluation function through the machine learning paradigm as the nature of the problem allows to generate large data to train

on using self-play methods. Recently, deep learning (LeCun et al., 2015; Schmid-huber, 2015) has proven a great fit for such a task as they provide excellent empirical generalization.

## 8.5   Lookahead Tree Size

With an exact value function, the size of the lookahead tree is irrelevant as even a single one-step lookahead is guaranteed to produce an optimal policy. But once we have only approximate values, one would expect larger trees to result in better policy. Indeed in practice, the size of the tree does matter. While the are good reasons for this to be the case, there is no guarantee for this behavior (Nau, 1982).

Furthermore, size is not the only parameter that matters in practice and structure is often also critical. We often want to search deeper where the current policy is likely to play or where we are uncertain about the values (similarly to how humans perform search in games). For example in chess, professional players search very deep where they believe the situation to be interesting. It is thus also common to dynamically expand the lookahead tree during the search rather than to construct fixed-sized lookaheads.

# 9. Examples

We now list some important milestones, where the combination of search and value functions led to historical achievements in perfect information games. While the search methods and value functions evolved over the years, the high level ideas can be traced to the very dawn of games research.

## 9.1 Samuel's Checkers

Samuel's checkers is an exciting application that dates all the way back to 1959 (Samuel, 1959). It is amazing to see how many of the ideas considered as new and modern already appeared in this program. While the strength of the program was relatively weak, one must realize how limited both hardware and software were at that time. And while search algorithms for perfect information games have made a great deal of progress both in strength and generality since then, the core principles of search combined with self-play learned value functions are already present here.

**Search Method** The search algorithm was a variant of minmax with Alpha-Beta pruning (Section 6.1.2). The lookahead tree was dynamically grown using rules based on the depth as well as some game specific properties.

**Value Functions** Perhaps the most impressive part were the value functions. The value was estimated via a polynomial combinations of hand-crafted game specific features. Not only were the weights adjusted through self-play — the learning method even used bootstrapping[1] techniques. Finally, the paper also describes a method of learning from prior human games.

## 9.2 Deep Blue

Being one of the most popular board game in the world, chess playing algorithms have always been of great interest to computer scientists. Shannon's work from 1950 already describes a simple variant of minimax search and even includes a section on approximate value function computed as a combination of game specific features (Shannon, 1950).

Almost 50 years after that — in 1997 — Deep Blue defeated the world champion Garry Kasparov (3.5–2.5). For the first time in history, machines defeated the human world champion in a regulation match (Campbell et al., 2002). In 1996, the previous version of Deep Blue lost to Kasparov by 2-4, making the second match particularly interesting.

**Search Method** Deep Blue used a massively parallel search that combined both hardware and software search. The search was a variant of AlphaBeta search with dynamically grown tree designed to be highly non-uniform. The

---

[1]The update/target for the value function being based on the value function estimate of the future states, see e.g. Sutton and Barto (2018).

software/hardware search allowed it to search up to 330 million chess positions per second.

**Value functions**  Deep Blue's evaluation was implemented in hardware. This made the evaluation particularly fast, but hard to modify. The evaluation function consisted of about $8,000$ different chess specific features ("passed pawns", "bishop pair") that were combined to produce the final value estimate. Making such a complicated game specific value function is a hard and time-consuming task. Authors themselves stated that "...we spent the vast majority of our time between the two matches designing, testing, and tuning the new evaluation function.".

## 9.3   TD-Gammon

TD-Gammon showed how powerful a combination of self-play learning and neural nets as value function approximation can be (Tesauro, 1995). It was the first successful application of this combination for large games, achieving impressive performance.

**Search Tree**  The search tree was a fairly small 2-ply lookahead, with larger trees suggested as a possible future improvement.

**Value functions**  TD-Gammon used a feed-forward fully conected neural network for its value function. The feature representation initially encoded raw board representation, and further game-specific hand-coded features were later added to improve the performance.

## 9.4   Alpha Go

AlphaGo further showed the power of combining deep learning and self-play. Go is a prime example of game where hand-crafting value functions or features is particularly difficult, and thus state of the art programs relied on rollout simulations to estimate the state values.

But the strength of such agents was far below professional players, and many predicted that it would take further decades for computers to match the top humans. AlphaGo shocked many by reaching this milestone in 2016 (Silver et al., 2016).

**Search Method**  The search used is a variant of MCTS with a p-UCB formula to guide the simulations.

**Value functions**  Value functions are represented with deep convolutional neural network (LeCun et al., 2015; Schmidhuber, 2015), and were pre-trained on human data through supervised training. Self-play was then used for further improvement. While some game specific features were still present in this version, further improvements made the agent both more general and stronger.

## 9.5   From AlphaGo to MuZero

Even though AlphaGo was a huge achievement, there were still some limitations. It could only play go, it needed human data for the initial supervised training and the value function represented by a neural network received some hand-crafted go features. Later work removed all of these limitations. AlphaGoZero uses no human data and removes hand-crafted features, while achieving even stronger performance (Silver et al., 2017b). AlphaZero then achieves master level perfomance on go, chess, and shogi — using a single algorithm (Silver et al., 2017a). This work has also now been reproduced by Leela Zero, one of the strongest and opensourced chess engines (Pascutto, 2019).

The latest representative of this lineage, MuZero makes the approach even more general, as it does not need the rules of the game to construct a game tree for search. The rules are learned during the self-play and the search uses only the learned model during the planning. This means that MuZero does not need to access the environment during its search, and only does so when it is taking its action (Schrittwieser et al., 2020).

# 10. Summary

In the first part of this book — perfect information games — we build a framework that allows us to build a particular view on search in games. First, we introduced optimal policies and their properties making them desirable to be followed. We then presented offline algorithms producing optimal policies, i.e. algorithms that solve games. Next, we described how sound online search algorithms follow these optimal offline strategies.

In the second part of this book — imperfect information games — we will build upon the formalism, intuition and framework. We will see that a combination of search and value functions is indeed possible and is also a powerful method in imperfect information. But we will also see that sound search is substantially harder in the imperfect information case.

# Part II

# Imperfect Information

# 11. Introduction

In the second part of this thesis, we will revisit the concepts introduced in perfect information settings. We will see that these concepts need to be generalized in order to be either well-defined or to produce the desired policies in imperfect information settings. But after we do so, we will indeed be able to construct sound search methods for imperfect information games, resulting in agents that can outplay even the best human players.

## 11.1 Structure

We start again with a quick overview of the chapters.

**Example Games**    Chapter 12 introduces the games used in the following chapters.

**Formalisms**    Chapter 13 includes formal models of imperfect information environments. Unlike in perfect information, we present multiple models and discuss their connections. Namely, we introduce factored-observation stochastic games — one of the contributions of this thesis.

**Offline Policies**    Chapter 14 argues that the same worst-case reasoning holds in imperfect information. The concepts of minmax and Nash equilibrium are still guaranteed to exist, with all the same desirable properties as in perfect information. But unlike in perfect information, the policies might have to be stochastic. We present some examples as to why the randomization is necessary, and introduce the link between the level of uncertainty in the game and the amount of randomization.

**Sub-games**    Chapter 15 generalizes the notion of sub-games. Unlike in perfect information, a single state is not sufficient to construct a well-defined sub-problem to be reasoned about. A distribution over the consistent states of both players is required as the agents can have different observations. We show how the notion of public information makes the construction of these consistent states particularly simple. Finally, we also generalize the connected concept of value functions to imperfect information sub-games.

**Offline Solving I**    Chapter 16 presents multiple offline algorithms producing optimal policies. These include direct optimization against the best-responding opponent, self-play via independent reinforcement learning, and fictitious play.

**Offline Solving II - Regret Minimization**    Chapter 17 deals with the regret minimization framework and its connections to optimal policies. Importantly, we introduce the powerful idea of counterfactual regret minimization that allows the use of regret minimization in sequential decision making by decomposing the regret to partial regrets in individual information states. Counterfactual regret

minimization will be used as key building blocks for the online search settings in the following chapters. We also include the family of sampling Monte Carlo variants and the latest addition to this family — methods that speed up the convergence by lowering the per-iteration variance. These low-variance methods are another contribution of this thesis.

**Offline Solve III - Approximate and Abstraction Methods**  Chapter 18 discusses methods that do not employ the idea of online search, but still allow one to deal with large games where tabular methods are not feasible.

**Online Settings**  Chapter 19 revisits the online setting, with some minor modifications to the required definitions to match the formalism of imperfect information. It includes some interesting counter-examples that are only possible in imperfect information.

**Search**  Chapter 20 is arguably the key chapter of this thesis, the culmination of the building blocks presented up to this point. First, we show that many properties we depended on in perfect information settings simply break in imperfect information. We then present the building blocks of re-solving and gadget, including another contribution of this thesis — safe refinement of sub-games. Step by step, we then put the building blocks carefully together into the final algorithm — continual resolving. Continual resolving allows for sound search in imperfect information, and was a crucial contribution of DeepStack. Chapter 20 and DeepStack are the main contributions of this thesis.

**DeepStack**  Chapter 21 presents DeepStack — the first computer agent to introduce the combination of sound search and learned value functions for poker, beating professional human players in no-limit Texas hold'em poker. As human evaluation in no-limit poker is inherently noisy, we also describe some techniques that allow for provably unbiased variance reduction. The latest of such techniques is AIVAT — our final contribution that resulted in an impressive 85% reduction in standard deviation when evaluating the DeepStack results. We finish with a list of some of the latest search algorithms.

# 12. Example Games

This chapter briefly introduces some imperfect information games used in the second part of the book. The quintessential example is poker, where we do not get to see opponent's cards. Consider the groundbreaking publication of Oskar Morgenstern and John von Neumann "Theory of Games and Economic Behavior" (Morgenstern and Von Neumann, 1953). This book laid the ground of modern-day game theory, and contains a full section (section 19, over 30 pages) dedicated to Poker.

The importance of poker resulted in substantial body of research where poker games were the only domain used for evaluating the algorithms. To emphasize the generality of the algorithms presented here, we decided to also include a game with structure drastically different to poker.

## 12.1   Kuhn Poker

Kuhn poker is a very simple poker game with only 3 cards in the deck and a single betting round (Kuhn, 1950). The deck includes only three cards: (J)ack, (Q)ueen and (K)ing. Both players are dealt a single card out of a deck, they each put one chip the table (ante) and a simple betting round follows. First player either (b)ets a chip (adding one more chip on the table) or (c)hecks. If the first player bets, the second player can either (c)all or (f)old — both actions terminating the game. If the first player checks, the second player can either bet one chip or check. Check terminates the game, while for bet, the first player acts again and can either check or fold — both actions terminating the game. If any player folded, they lose the chips on the table. Otherwise, the player holding the higher card wins these chips. We will present figures of small poker games as soon as we introduce the necessary formalism in Section 13.3.7 and 13.3.7.

## 12.2   Leduc Poker

Leduc poker is a slightly larger poker game often used as a test domain in imperfect information games. The card deck consists of 6 cards (two suits and three ranks), and there are two (limited) betting rounds with a single public board card (Southey et al., 2005). Note that while the game is larger than Kuhn poker, it is still very small.

## 12.3   Limit nad No-limit Texas Hold'em Poker

Texas hold'em poker is played with the standard 52 card deck. The game progresses in four betting rounds: i) pre-flop ii) flop iii) turn and iv) river. Cards are dealt at the beginning of each round. In the pre-flop, each player is dealt two private cards (hand). In the later rounds, public (facing up and observable to all players) cards are dealt on the table: three on the flop, one on the turn and one on the river. During each betting round, the players alternate taking one of the three action types: i) fold ii) check/call and iii) bet/raise. By folding, the player

gives up and loses the money wagered up to this point (pot). Calling matches the bet of the opponent, and check is the initial call action when the player is facing no bet/raise. Raise increases a player's wager, and bet is the initial raise when there is no opponent's bet to be raised. The round ends when both players match their bet, and the game terminates when either of the player folds or when the last round (river) ends. When neither player folds before the end of the river, both players reveal their hand and the player with the strongest combination of the private and public cards wins the pot. For the ranking of the possible card combinations and more detailed rules see e.g. Harroch and Krieger (2007).

Finally, the difference between limit and no-limit poker is that in limit poker, there is a single bet size the player can make. In no-limit, the player is free to make a bet up to their stack size, and all-in bet then refers to agent betting their entire stack

## 12.4  Graph Chase

Graph chase games are a simple, security-like games. Players get to move their stones from node to node, alternating their moves. One player (evader) controls a single invisible stone and is trying to escape the opponent. The second player (chaser) controls multiple stones (the location of these stones is publicly observable) and is trying to catch the invisible stone (by moving a stone on to the location of the invisible stone). The evader wins if they gets to survive a specified number of moves, and loses otherwise. While the rules can surely be made more complicated, this simple variant still makes an interesting game.

Figure 12.1 shows a small instantiation of this game, with two stones of the chaser and a small graph. We refer to this game as the Glasses due to the shape of the graph.

### 12.4.1  Contrast to Poker

There are some important conceptual difference to poker. In poker, all actions are publicly observable, making some of the decomposition constructions easier. In graph chase games, the chaser does not get to see the exact action (edge) the evader took. Furthermore, the number of states the player can be in varies as the game progresses. In poker, this corresponds to the possible hands of a player, which is fixed in individual betting rounds. These distinctions are important as we need our algorithms to handle this general case.
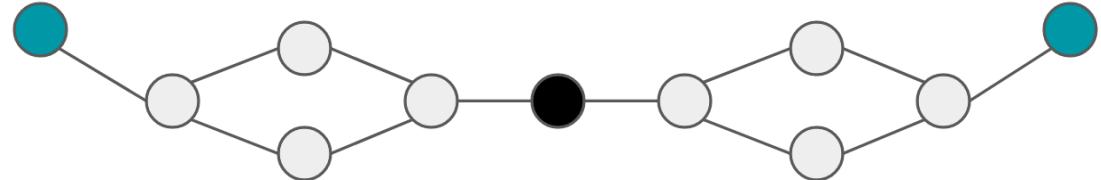


Figure 12.1: Glasses: graph chasing game on a graph with only 11 nodes. Evader starts in the middle (black stone) while the chaser has two (blue) stones at the edges. The evader wins if they get to survive for 3 turns.

# 13. Formalisms

In this chapter, we introduce three distinct formal models of imperfect information games. The first one — matrix games — is a simple yet powerful model for simultaneous decision making, where the players get to act at the same time. The second model — extensive form games - is the most popular formal model of sequential decision making in imperfect information. It is a natural generalization of the game tree from perfect information setting. Unfortunately, some design choices of this model make it ill-suited for formal analysis of the search algorithms, motivating the next model.

The third model — Factored Observation Stochastic Games (FOSG) — is a recent model build with the needs of search in mind and also one of the contributions of the thesis. Just like extensive form games, it models sequential decision making in imperfect information. The defining properties that make it suitable for modern search methods are i) factoring observations to private and public components, and ii) the notion of an agent's state even if they are not to act.

Note that any sub-problem decomposition and search only really makes sense in sequential decision making. Thus while we use matrix games to introduce some of the essential theory, sub-problems and search methods will be only discussed for FOSGs.

## 13.1 Matrix Games

It turns out that the simplest formalism for imperfect information games is based on simultaneous moves. That is, both players get to make their action at the same time. At first, it might not be clear why this has anything to do with imperfect information. The reason is simple — we can think of such interactions as one of the players acting first, while hiding the action from the opponent. The opponent then gets to act without knowing what the first player did.

Matrix games (also referred to as normal form games) is a simple formalism, where the first player (the row player) chooses a row, while the second player (the column player) chooses a column (Definition 16). The reward of the the first player is then computed as $xXy^\top$.

**Definition 16.** *Matrix game consists of:*

- $\mathcal{N} = \{1, 2\}$ *is the set of players.*

- $\mathcal{A}_1, \mathcal{A}_2$ *are the player's set of legal actions.*

- *Player strategies* $x \in \Delta(\mathcal{A}_1)$, $y \in \Delta(\mathcal{A}_2)$.

- *The reward matrix* $\underset{|\mathcal{A}_1| \times |\mathcal{A}_2|}{X}$.

### 13.1.1 Example

Table 13.1 includes the matrix game for the rock-paper-scissors game. In this case, both row and column player have the same set of actions, $A_1 = A_2 = \{R, P, S\}$. As a sanity check, the row's player utility under the following strategy profile $(\pi_1 = (0.2, 0.2, 0.6), \pi_2 = (0.4, 0.2, 0.4))$ is $\pi_1 A \pi_2' = 0.08$.

|   | R | P | S |
|---|---|---|---|
| R | 0 | 1 | -1 |
| R | -1 | 0 | 1 |
| R | 1 | -1 | 0 |

Table 13.1: Matrix game representation of the (R)ock-(P)aper-(S)cissors game.

## 13.2 Extensive Form Games

Extensive form games are a natural extension of game trees to imperfect information environments, allowing for sequential interactions. In imperfect information, players do not directly observe the exact state of the game and thus there might be mutliple world states the players can not tell apart. Extensive form games achieve this by grouping such states into so-called information sets (or states). The agent's policy is then defined within these information states.

Extensive form games date all the way back to Von Neumann (Morgenstern and Von Neumann, 1953), and the model has been widely adopted by the game theory community. The importance of the original publication is illustrated by the fact that the commemorative edition was published more than 50 years after the original version (Von Neumann and Morgenstern, 2007). Interestingly enough, the book contains a full section (section 19, over 30 pages) dedicated to Poker. And in 2015, a limit Texas hold'em poker become the largest imperfect information game played by humans to be solved (Bowling et al., 2015) - modeled by this very formalism.

**Definition 17.** *(Osborne and Rubinstein, 1994, Definition 200.1) An extensive form game is a tuple $(\mathcal{H}, \mathcal{Z}, \mathcal{A}, \mathcal{N}, p, \pi_c, u, \mathcal{I})$, where:*

- *$\mathcal{H}$ is the set of histories, representing sequences of actions.*

- *$\mathcal{Z}$ is the set of terminal histories (those $z \in \mathcal{H}$ which are not a prefix of any other history). We use $g \sqsubseteq h$ to denote the fact that $g$ is equal to or a prefix of $h$.*

- *For a non-terminal history $h \in \mathcal{H} \setminus \mathcal{Z}$, $\mathcal{A}(h) := \{a \mid ha \in \mathcal{H}\}$ is the set of actions available at $h$.*

- *$\mathcal{N} = \{1, \ldots, N\}$ is the player set. In addition, $c$ is a special player, called "chance" or "nature".*

- *$p : \mathcal{H} \setminus \mathcal{Z} \to \mathcal{N} \cup \{c\}$ is the player function partitioning non-terminal histories into $\mathcal{H}_i$, depending on which player acts at $h$.*

- *The strategy of chance is a fixed probability distribution $\pi_c$ over actions in $\mathcal{H}_c$, $\pi_c(h) \in \Delta(\mathcal{A}(h))$.*

- *The utility function $u = (u_i)_{i \in \mathcal{N}}$ assigns to each terminal history $z$ a reward $u_i(z) \in \mathbb{R}$ received by player $i$ upon reaching $z$.*

- *The information-partition $\mathcal{I} = (\mathcal{I}_i)_{i \in \mathcal{N}}$ captures the imperfect information of $G$. For each player, $\mathcal{I}_i$ is a partition of $\mathcal{H}_i$. If $g, h \in \mathcal{H}_i$ belong to the*

*same $I \in \mathcal{I}_i$ then $i$ cannot distinguish between them. For each $I \in \mathcal{I}_i$, the available actions $\mathcal{A}(h)$ must the same for each $h \in I$, and we overload $\mathcal{A}(\cdot)$ as $\mathcal{A}(I) = \mathcal{A}(h)$.*

### 13.2.1 Strategies

Behavioral strategies prescribe behavior in all the individual information sets (states). It is a mapping from an information state to a distribution over the actions $\pi_i : I \in \mathcal{I}_i \to \Delta A(I)$. Section 16.3.1 will present another option for strategy representation — the sequence form.

### 13.2.2 Example



Figure 13.1: Extensive form game for Rock-Paper-Scissors. Individual states (histories) are labeled with the corresponding history (there are 13 histories in total, 9 of which are terminal). Below the terminal states, we also include the player 1 terminal utility. Finally, the dotted lines connect histories into information sets/states. In this case, both players have a single information set.
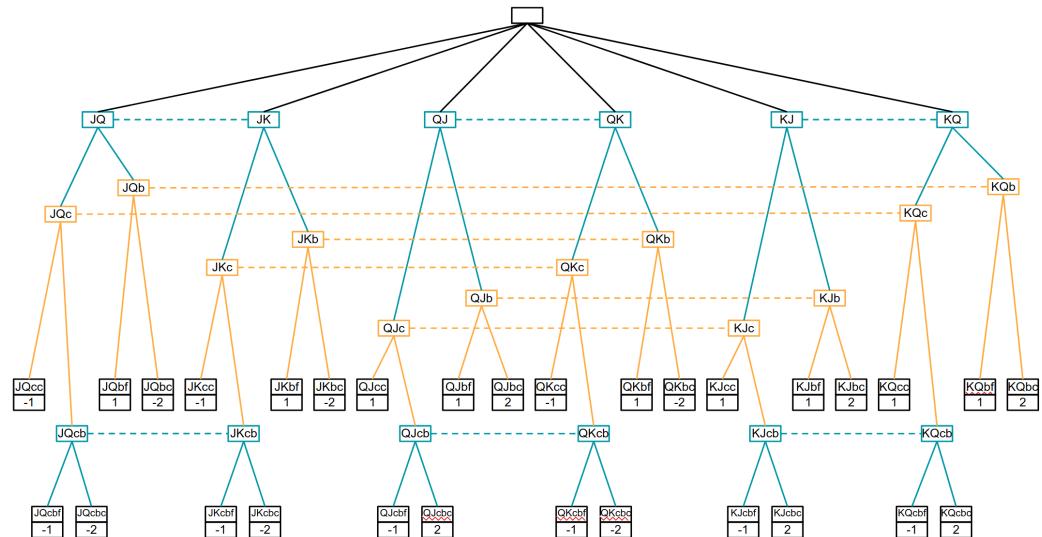


Figure 13.2: Extensive form game for the Kuhn Poker. Dashed lines connect histories grouped within an information set.

The first example illustrates how to capture the simultaneous move game of rock-paper-scissors in the inherently sequential model of extensive form games (Figure 13.1). While illustrated on rock-paper-scissors, it is easy to imagine the

general construction for an arbitrary matrix game (Corollary 12.1). The idea is to simply let one player to act first, group all the subsequent histories into a single information set so that the opponent does not know what action the first player made (Osborne and Rubinstein, 1994).

The second example is then the sequential game of Kuhn poker, illustrated in Figure 13.2.

**Corollary 12.1.** *Any matrix game can be converted to an equivalent extensive form game having (up to a constant factor) the same size.*

### 13.2.3 Extensive Form Games to Matrix Games

We have already seen how to convert any matrix game to an extensive form game (Corollary 12.1). Lemma 13 then shows that the opposite conversion is also possible under the perfect recall (Section 13.3.1) assumption. Unfortunately, the resulting matrix game can be exponentially larger than the original sequential representation of extensive form games. See Figure 13.3 for an example transformation (taken from Schmid (2013)).

**Lemma 13.** *(Osborne and Rubinstein, 1994) Given any two-player extensive form game with perfect recall, it's possible to create an equivalent normal form game.*
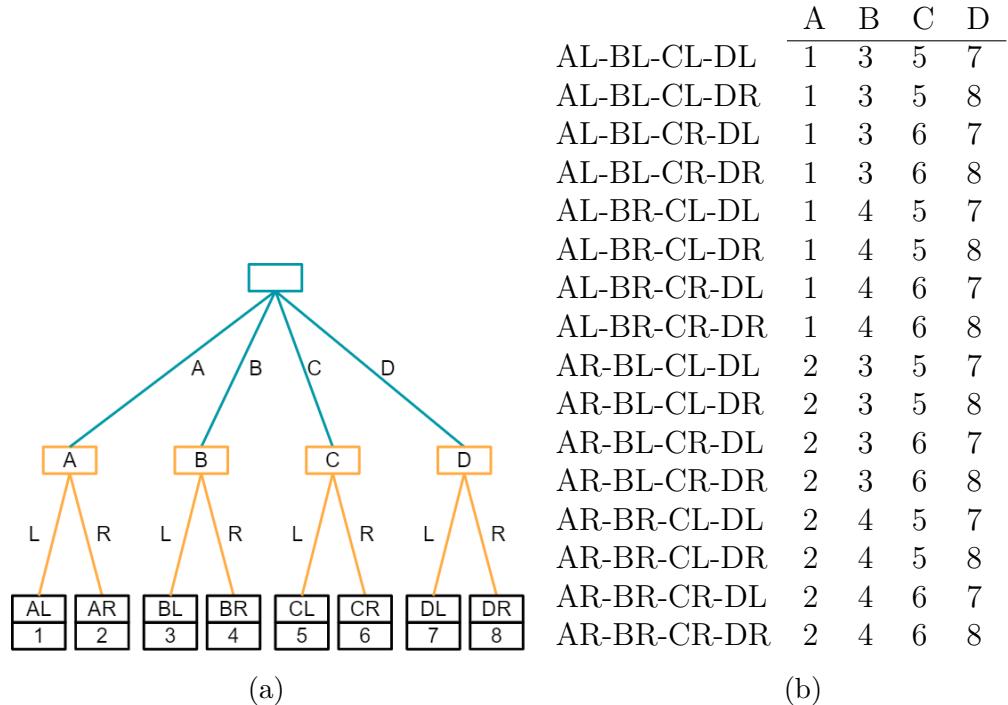


|           | A | B | C | D |
|-----------|---|---|---|---|
| AL-BL-CL-DL | 1 | 3 | 5 | 7 |
| AL-BL-CL-DR | 1 | 3 | 5 | 8 |
| AL-BL-CR-DL | 1 | 3 | 6 | 7 |
| AL-BL-CR-DR | 1 | 3 | 6 | 8 |
| AL-BR-CL-DL | 1 | 4 | 5 | 7 |
| AL-BR-CL-DR | 1 | 4 | 5 | 8 |
| AL-BR-CR-DL | 1 | 4 | 6 | 7 |
| AL-BR-CR-DR | 1 | 4 | 6 | 8 |
| AR-BL-CL-DL | 2 | 3 | 5 | 7 |
| AR-BL-CL-DR | 2 | 3 | 5 | 8 |
| AR-BL-CR-DL | 2 | 3 | 6 | 7 |
| AR-BL-CR-DR | 2 | 3 | 6 | 8 |
| AR-BR-CL-DL | 2 | 4 | 5 | 7 |
| AR-BR-CL-DR | 2 | 4 | 5 | 8 |
| AR-BR-CR-DL | 2 | 4 | 6 | 7 |
| AR-BR-CR-DR | 2 | 4 | 6 | 8 |

(a)                                             (b)

Figure 13.3: (a) Extensive form game game. (b) Corresponding matrix game.

## 13.3 Factored-Observation Stochastic Games

Factored observation stochastic games (FOSG) is a recently introduced formalism, and also one of the contributions of the thesis (Kovařík et al., 2019). We first

discuss the limitations of extensive form games and motivate the new formalism that builds on the notion of observations. We then formally introduce FOSG and finish with connections between the FOSG and EFG.

### 13.3.1 Limitations of EFG Formalism

The defining distinction of EFGs compared to perfect information game trees is the information partition forming information sets (states). Information states group histories (world states) that a player cannot tell apart. While this allows for a very general class of games, it turns out that grouping arbitrarily states of the acting player is both too general and loses important information.

It is too general as it allows for imperfect recall, which forces the agent to forget previously known information (e.g. past actions taken). This happens when information states group history with its prefix (i.e. grouping a node with a previous one). Imperfect recall is troubling as it not only is unrealistic, it has unfortunate complexity consequences. Many basic concepts including best response become complicated (Piccione et al., 1996; Piccione and Rubinstein, 1997), and even the existence of Nash equilibria becomes NP-hard (Hansen et al., 2007). It is thus common to restrict the set of games to only perfect recall. Recently, non-timeability has been suggested as another unrealistic property allowed by this model (Jakobsen et al., 2016). The observation is that since the passage of time is observable to agents, some extensive form games can not be implemented. A game is said to be timeable if it has an exact deterministic timing and 1-timeable if each label is exactly one higher than its parent's label (Definition 18). Similarly to perfect recall, authors argue that timeability should be a common assumption.

EFGs lose crucial information inherently present in many environments — the notion of observations, specifying what and when they were observed by the agents. Information states represent information available to the acting player, but there is no notion of information state for the non-acting player. This is problematic as when we are to play, we might need to reason about the states the opponent can be in. Furthermore, EFGs make no explicit distinction between private and public information. But the knowledge about what information is available to the non-acting player as well as the factoring of information into private and public are critical concepts for sound search method. While it is common to recover those concepts in specific cases (e.g. if we are building a poker agent) (Burch et al., 2014), it turns out that it is impossible to do properly in general (Kovařík and Lisý, 2019).

**Definition 18.** *(Jakobsen et al., 2016, Timeability) For an extensive-form game, a deterministic timing is a labelling of the nodes with non-negative real numbers such that the label of any node is at least one higher than the label of its parent. A deterministic timing is exact if any two nodes in the same information set have the same label.*

### 13.3.2 Augmented EFG

Recovering the information lost by the EFG formalism essentially led to the introduction of augmented information sets and public tree (Johanson et al.,

2011; Burch et al., 2014). These additional concepts then result in the augmented extensive form games (Definition 19).

**Definition 19.** *(Kovařík et al., 2019, Definition 3.5) Augmented extensive form game is a tuple $(\mathcal{H}, \mathcal{Z}, \mathcal{A}, \mathcal{N}, p, \pi_c, u, \mathcal{I})$ where all objects are as in extensive form game except for the information set partitioning $\mathcal{I}$:*

- *$\mathcal{I} = (\mathcal{I}_1, \mathcal{I}_2, \mathcal{I}_{pub})$ is a collection of partitions of $\mathcal{H}$ where each $I_i$ is a refinement of $\mathcal{I}_{pub}$*

Note that there are two distinctions to the standard information sets. First, the information sets partition $\mathcal{H}$ rather than just $\mathcal{H}_i$ and information sets are thus defined even when the player is not acting. Second, $\mathcal{I}_{pub}$ defines public states and a corresponding public tree where the individual player's partitions are a refinement of this public partitioning.

But this recovery is only possible in specific cases, using game-specific concepts (e.g. dealing cards in poker). Furthermore, on top of the usual restrictions of perfect recall and timeability, additional restriction is then assumed. Augmented EFG is said to not have thick public sets if no element of $\mathcal{I}_{pub}$ (and hence of $\mathcal{I}_i$) contains both some $h$ and $h'$ for $h' \sqsubset h$. This additional technicality is essentially the perfect recall analog for augmented information sets.

### 13.3.3 FOSG Definition

Rather than using complex combination of restrictions and augmentations of extensive form games, we argue that a simple observation-based model naturally describes the domain and preserves the necessary information. Motivated by the issues of extensive form games (Section 13.3.1), factored observation stochastic games is a recently proposed formalism of multi-agent imperfect information environments (Kovařík et al., 2019).

The formalism is a variant of partially observable stochastic games (Hansen et al., 2004a), where the game consists of underlying world states, and the probabilistic transition to a next world state is a function of the actions taken by all the agents. As the game moves from one world state to another, players do not get to directly observe the underlying state. Rather, they receive an observation that is factored to a private observation and a public observation. The notions of transitions, world states and observations also make this formalism more familiar to the the reinforcement learning community.

**Definition 20.** *(Kovařík et al., 2019) Factored Observation Stochastic Game is a tuple $G = (\mathcal{N}, \mathcal{W}, w^o, \mathcal{A}, \mathcal{T}, \mathcal{R}, \mathcal{O}, \mathbb{O})$*

- *$\mathcal{N} = \{1, \dots, N\}$ is the player set.*

- *$\mathcal{W}$ is the set of world states and $w^0 \in \mathcal{W}$ is a designated initial state.*

- *$\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_N$ is the space of joint actions.*

  - *The subsets $\mathcal{A}_i(w) \subset \mathcal{A}_i$ and $\mathcal{A}(w) = \mathcal{A}_1(w) \times \dots \times \mathcal{A}_N(w)$ specify the (joint) actions legal at $w$.*

  - *For $a \in \mathcal{A}$, we write $a = (a_1, \dots, a_N)$.*

- $\mathcal{A}_i(w)$ *for* $i \in \mathcal{N}$ *are either all non-empty or all empty. A world state with no legal actions is terminal.*

- *After taking a (legal) joint action* $a$ *at* $w$*, the transition function* $\mathcal{T}$ *determines the next* $w' \sim \mathcal{T}(w, a) \in \Delta(\mathcal{W})$.

- *We write* $R = (R_1, \ldots, R_N)$*, where* $R_i(w, a)$ *is the reward* $i$ *receives when a joint action* $a$ *is taken at* $w$.

- $\mathcal{O} = (\mathcal{O}_{\mathrm{priv}(1)}, \ldots, \mathcal{O}_{\mathrm{priv}(N)}, \mathcal{O}_{\mathrm{pub}})$ *is the observation function, and* $\mathbb{O} = (\mathbb{O}_{priv(1)}, \mathbb{O}_{priv(2)}, \mathbb{O}_{pub})$ *are the observation sets. The observation function* $\mathcal{O}_{(\cdot)} : \mathcal{W} \times \mathcal{A} \times \mathcal{W} \to \mathbb{O}_{(\cdot)}$ *specifies the private observation that* $i$ *receives, resp. the public observation that everybody receives, upon transitioning from world state* $w$ *to* $w'$ *via some* $a$.

  - *For each* $i$*, we write* $\mathcal{O}_i(w, a, w') = \left( \mathcal{O}_{\mathrm{priv}(i)}(w, a, w'), \mathcal{O}_{\mathrm{pub}}(w, a, w') \right) \in \mathbb{O}_i := \mathbb{O}_{\mathrm{priv}(i)} \times \mathbb{O}_{\mathrm{pub}}$.

  - *We assume that at the start of the game, each player receives some* $\mathcal{O}_i(w^0)$.

The game then starts in the initial state $w^0$ and follows in turns. In each turn, each player $i$ player select an action $a_i \in \mathcal{A}_i(w)$, resulting a joint action $a = (a_i)_{i \in \mathcal{N}}$. The game then transitions to a new state $w'$ $\mathcal{T}(w, a)$. This transition generates an observation $\mathcal{O}(w, a, w')$, from which each player receives $\mathcal{O}_i(w, a, w') = (\mathcal{O}_{\mathrm{priv}(i)}(w, a, w'), \mathcal{O}_{\mathrm{pub}}(w, a, w'))$ (i.e., the public observation together with their private observation). Finally, each player is assigned the reward $\mathcal{R}_i(w, a)$. This process repeats until a terminal state is reached.

As an example, consider the rock-paper-scissors. The set of world states is $\mathcal{W} = \{0, r, p, s, rr, rp, rs, pr, pp, ps, sr, sp, ss\}$, with the initial state $w^0 = 0$. As only a single player acts at each time step, the non-acting player has a single dummy action $\mathcal{A}_2(r) = \{r, p, s\}, \mathcal{A}_1(r) = \{-\}$. The transition function is fully deterministic, e.g. $\mathcal{T}(r, (r, -)) = rr$. Finally, the observation sets in rock-paper-scissors are empty. In poker games, the private observations would encode the private cards, while the public observations provide information about the public (board) cards as well as about the actions taken (fold/call/bet) as the actions Texas hold'em poker are publicly observable.

### 13.3.4 Properties

As the agents receive observation during the game-play, the full action observation history is necessarily perfect recall. The formalism does not force agent into imperfect recall (and there are extensions allowing for imperfect recall).

As all the agents act all the time, there is no issue with ill-defined information (states) of the non-acting agents. This has additional benefit, as this implies a notion of "tick/time step" inferred by the number of actions taken, and thus non-timeable games are not possible.

Finally, the factorization to private information and public information allows for public states and public trees.

61

### 13.3.5 Connection to EFG

There are clear benefits of the observation based model of FOSGs — it is relatively simple, encodes the information otherwise lost by EFG, and it contains only timeable perfect-recall games. We now show that there is no loss of expressive power as any perfect recall timeable EFG can be represented.

First, observe that any FOSG corresponds to some augmented extensive form representation through the `FOGToAugEFG` transformation (the formal process can be found in (Kovařík et al., 2019)). Next transformation `AugEFGToEFG` simply forgets the $\mathcal{I}_{pub}$ and restricts $\mathcal{I}_i$ to $\mathcal{H}_i$. Finally, we define `FOGtoEFG = AugEFGToEFG ∘ FOGToAugEFG`. With these transformations in hand, we can state the crucial connections.

**Theorem 14.** *(Kovařík et al., 2019, Theorem 3.6) (FOSGs as augmented EFGs) Every FOSG G corresponds to an augmented EFG E = **FOGToAEFG**(G) with perfect-recall and no thick public sets. Moreover, any perfect-recall augmented EFG with no thick public sets can be obtained this way.*

**Theorem 15.** *(Kovařík et al., 2019, Theorem 3.10) (FOSGs are timeable EFGs) For any FOSG G, the game E = **AugEFGToEFG** ∘ **FOGToAugEFG** (G) is a 1-timeable perfect-recall EFG. Moreover, any 1-timeable perfect-recall EFG can be obtained this way.*

Finally, it is possible to show that 1-timeable EFG is no less general than any timeable game (Lemma 16).

**Lemma 16.** *(Kovařík et al., 2019, Lemma 3.9) Any classical timeable EFG can be made 1-timeable by adding chance nodes with a single noop action, where the resulting 1-timeable game is no more than quadratically larger.*

### 13.3.6 Derived Objects

We can now readily derive the trees of histories, cumulative rewards and information states, similar to the objects used in the extensive form games literature.

- History (trajectory) is a finite sequence $h = (w^0, a^0, w^1, a^1, \ldots, w^t)$, where $w^i \in \mathcal{W}$, $a^i \in \mathcal{A}(w^i)$, and $w^{i+1} \in \mathcal{W}$ is in the support[1] of $\mathcal{T}(w^i, a^i)$. We denote the set of all legal histories by $\mathcal{H}$[2], and the set of terminal histories as $\mathcal{Z}$.

  - Since the last state in each $h \in \mathcal{H}$ is uniquely defined, the notation for $\mathcal{W}$ can be overloaded to work with $\mathcal{H}$ (for example $\mathcal{A}(h) = \mathcal{A}(w^t)$).

  - We use $g \sqsubset h$ to denote the fact that $g$ is a prefix of $h$, and $g \sqsubseteq h$ to denote $g$ is a prefix of $h$ or equal to $h$.

- The cumulative reward (return or utility) of $i$ at $h$ is the reward accumulated up to this point $R_i(h) = \sum_{i=0}^t R_i(w^i, a^i)$.

---

[1]For finite $\mathcal{W}$, being in support of some probability measure $\mu \in \Delta(W)$ is equivalent to having a non-zero probability under $\mu$.

[2]To simplify the discussion, we assume that $\mathcal{H}$ is finite. This can always be enforced by using finite horizon $T$, but some domains satisfy this assumption naturally. Moreover, our results generalize into the standard $\gamma$-discounted rewards setting.

The combination of private and public observations induce the corresponding information states and trees, and we then use these to define the notion of a strategy.

- Player $i$'s action-observation history at $h$ is the sequence of the observations visible to that player (private and public) and the actions taken: $s_i(h) = (o_i^0, a_i^0, o_i^1, a_i^1, \ldots, o_i^t)$, where $o_i^k = \mathcal{O}_i(w^{k-1}, a^{k-1}, w^k)$. The space $\mathcal{S}_i$ of all such sequences can be viewed as the information state space of $i$.

  - We use $\mathcal{H}(s)$ to denote the set of compatible histories: $\mathcal{H}(s) = \{h \in \mathcal{H} | s_i(h) = s\}$
  - We assume that each $s_i \in \mathcal{S}_i$ determines which $i$'s actions $\mathcal{A}_i(s_i)$ are legal: $(\forall h \in \mathcal{H}_i(s_i)) : \mathcal{A}_i(h) = \mathcal{A}_i(s_i)$
  - We also use *sao* to denote a state after taking action $a$ in $s$ and observing $o$.

- A policy $\pi_i : s_i \in \mathcal{S}_i \mapsto \pi_i(s_i) \in \Delta(\mathcal{A}_i(s_i))$ is a mapping from an information state to the probability distribution over the actions in that state, and $\Pi_i$ is the set of all possible policies.

Finally, the factorization of each observation into the private and public parts allows us to define the tree of public states.

- The public state corresponding to $h = (w^0, a^0, w^1, a^1, \ldots, w^t)$ is the sequence $s_{\text{pub}}(h) := s_{\text{pub}}(s_i(h)) := \left(O_{\text{pub}}^0, O_{\text{pub}}^1, \ldots, O_{\text{pub}}^t\right)$ of all public observations corresponding to $h$.

  - The space $\mathcal{S}_{\text{pub}}$ of all public states is called the public tree.

It is often useful to model the stochasticity in the game as being caused by a chance player $c$ rather than via the stochastic transition function $\mathcal{T}$. The chance player receives observations and takes actions just like the other players, but its policy is fixed. The addition of $c$ causes no loss of generality since their actions can always be merged back into the transition function $\mathcal{T}$, but this viewpoint can be in some cases conceptually and formally more convenient. Finally, we also use $\Delta_R$ to denote the delta of maximim/minimum utilities in the game: $\Delta_R = \max_{z \in \mathcal{Z}} R_i(z) - \min_{z \in \mathcal{Z}} R_i(z)$.

### 13.3.7 State Notation and Views

Factored observation games naturally allow for different observer views. Each world state contains i) player 1 private information ii) player 2 private information iii) public information. Different agents then observe different subsets of that available information, corresponding

- Game state view: i) player 1 private information ii) player 2 private information iii) public information

- Player 1 view: i) player 1 private information iii) public information

| | Notation | $O_{pub}$ | $O_1$ | $O_2$ | Kuhn Example |
|---|---|---|---|---|---|
| $h \in \mathcal{H}$ | $[O_{pub}\|O_1\|O_2]$ | ✓ | ✓ | ✓ | $[b\|Q\|J]$ |
| $s \in \mathcal{S}_1$ | $[O_{pub}\|O_1\|]$ | ✓ | ✓ | | $[b\|Q\|]$ |
| $s \in \mathcal{S}_2$ | $[O_{pub}\|\|O_2]$ | ✓ | | ✓ | $[b\|\|J]$ |
| $s_{pub} \in \mathcal{S}_{pub}$ | $[O_{pub}\|\|]$ | ✓ | | | $[b\|\|]$ |

Table 13.2: Different states/views are defined by different subsets of the information/observation. Kuhn example shows a state where the first player is dealt (Q)ueen, the second player is dealt (J)ack and the first player made a publicly observable action (b)et (see also Figure 13.4).

| | Histories | Infostates | Public States |
|---|---|---|---|
| Kuhn poker | 64 | 20 | 12 |
| Leduc poker | $9,487$ | $1,398$ | 468 |
| Glasses | $53,907$ | $11,699$ | $6,559$ |
| Limit Texas Hold'em | $3 \cdot 10^{17}$ | $3 \cdot 10^{14}$ | $3 \cdot 10^{11}$ |
| No-limit Texas Hold'em | $6 \cdot 10^{164}$ | $6 \cdot 10^{161}$ | $6 \cdot 10^{158}$ |

Table 13.3: The no-limit poker variant is the one used at the Annual Computer Poker Competition (ACPC) since 2010 and its size was computed in (Johanson, 2013).

- Player 2 view: ii) player 2 private information iii) public information)

- Public observer view: iii) public information

Thanks to the perfect recall property, all these views imply a well-formed trees about how the game progresses from the perspective of the respective viewers. Furthermore, this allows for a convenient notation for identifying states. We will use $[O_{pub}\|O_1\|O_2]$, $[O_{pub}\|O_1\|]$, $[O_{pub}\|\|O_2]$, $[O_{pub}\|\|]$ to identify a concrete world-state, player 1 state, player 2 state and public state respectively (see Table 13.2).

### 13.3.8 Example

Figure 13.4 shows the different views for the game of Kuhn poker. See that as the game view, player's view and the public view contain decreasingly less information, we are "zooming in and out" of the respective states. Table 13.2 then shows state notation for a concrete Kuhn state and Table 13.3 reports number of world states, infostates and public states for the imperfect information games used in this thesis.

Figure 13.4: Different levels of information for the Kuhn poker. Top left: Public tree depicting the publicly observable information (i.e. the public view). Top right: Action-observation sequence tree of player 1 (i.e., player 1 view). Bottom: The world-state tree of the game (i.e. the game view). Red: Example states used in Table 13.2.

## 13.4 History of Formalisms

As discussed in Section 1.7, there are two main communities that study sequential decision-making in multi-agent environments — game theory and reinforcement learning. We include a brief history of the formalisms used by the communities.

Extensive Form Games are a natural extension of game trees to imperfect information environments, dating back to 1953 (Morgenstern and Von Neumann, 1953). Factored Observation Stochastic Games are a modern formalism that reflects the lessons learned from working with extensive form games as well as having search methods in mind. The formalism is based on the stochastic games, and aims to bring the communities closer as the next research frontiers require insights and collaboration of both communities.

Reinforcement learning formalisms can be traced back to Markov decision processes (MDP) (Bellman, 1957) and stochastic games (Shapley, 1953). While the MDP formalism describes single-agent perfect information domains, researchers soon generalized MDPs into partially observable POMDPs (Astrom, 1965). Kaelbling et al. (1998) then used POMDPs to present novel algorithms for planning and acting, using value functions and sub-problems as crucial concepts. While stochastic games (SGs) were introduced prior to MDPs, they can be well thought of as "multiagent MDPs" (SG is a FOSG with empty observation sets $\mathbb{O}_{priv(i)} = \mathbb{O}_{pub} = \emptyset$). Indeed, this connection allowed (Littman, 1994) to generalize the Q-learning algorithm from MDPs to stochastic games. When it comes to combining multiple agents and imperfect information, researchers usually use partially observable stochastic games (Hansen et al., 2004b) (FOSG with no notion of public observations $\mathbb{O}_{pub} = \emptyset$)) and Decentralized POMDPs (POMDP

where the agents' goal is shared: $R_i = R_j \ \forall i, j \in \mathcal{N}$) (Bernstein et al., 2002).

# 14. Offline Policies

As the imperfect information games formalisms are strictly more general then their perfect information counterparts, we will revisit the concepts of optimal policies introduced in Chapter 4 and see if and how these concepts transfer to imperfect information.

Do the optimality results from perfect information settings hold in imperfect information? As a reminder, the settings we are dealing with are two-player, zero sum, perfect recall and finite-horizon games.

- Is a maximin policy still guaranteed to exist?

- Is a Nash equilibrium policy still guaranteed to exist?

- Do maximin and Nash equilibria collapse?

- Are there still the same worst-case garantees (as in Corollary 5.2)?

- Are the optimal policies still deterministic?

## 14.1   Best Response

The definition of best response in imperfect information games remains unchanged the perfect information (Section 4.1).

**Definition 4** (Best Response)**.** *Best response against a policy $\pi_i$ is:*

$$\arg\max_{\pi_{-i} \in \Pi_{-i}} R_{-i}(\pi_i, \pi_{-i})$$

And just as in perfect information, best response can always be a deterministic policy (Lemma 17). This follows from the fact that fixing the opponent turns the game into a perfect information single-agent environment (Bowling, 2003). To compute the best response, we use the fixed policy of the opponent to compute the transition probabilities of the agent's infoset tree. Using those transition probabilities, we then simply traverse the tree bottom-up and greedily, selecting the action with the highest expected utility (Algorithm 4).

**Lemma 17.** *There is always a deterministic best response.*

### 14.1.1   Distinction to Perfect Information

There is a very important distinction to perfect information in the way the best response is calculated. Previously, we needed to consider the opponent's strategy only in all the future states. This is because this provided a sufficient information about the transition probabilities between these future states. There was no need to consider how the opponent plays prior to a state and thus a truly bottom-up pass was possible in a perfect information best response (Algorithm 1).

This is not the case in imperfect information. The issue is that the state-action-state transition dynamic depends on the distribution over the individual

Figure 14.1: Best response computation for the second player in rock-paper-scissors. Single infostate of the player consists of three different histories, and the utility of the policy depends on distribution over these histories. That in turn depends on the strategy of the opponent in previous states.

histories grouped within the infostate. That distribution in turn depends on the past behavior of the opponent. Consider rock-paper-scissors from the perspective of the second player (i.e. computing best response against a fixed strategy of the first player). To compute the probability of reaching the terminal state with utility 1 after the action rock, we need the distribution over the three possible histories grouped in the infostate. And these three histories correspond to the first player taking rock, paper or scissors (Figure 14.1).

See that the Algorithm 4 indeed first needs to compute the reach probabilities for all the histories, propagating them down the tree. We will see important consequences of this property on the imperfect information sub-games in Section 15.

## 14.2 Maximin

Same as in perfect information, the maximin strategy maximizes the agent's value against the worst case opponent.

**Definition 5** (Maximin Policy)**.** *Maximin policy of a player i is:*

$$\arg\max_{\pi_i \in \Pi_i} \min_{\pi_{-i} \in \Pi_{-i}} R_i(\pi_i, \pi_{-i}) = \arg\max_{\pi_i \in \Pi_i} BRV_i(\pi_i) \qquad (4.1)$$

## 14.3 Nash Equilibrium

The concept and definition of Nash equilibria also remains unchanged.

**Definition 7** (Nash Equilibrium)**.** *Strategy profile $(\pi_i, \pi_{-i})$ forms a Nash equilibrium if none of the players benefit by deviating from their policy.*

$$\forall i \in N, \forall \pi_i' \,:\, R_i(\pi_i, \pi_{-i}) \geq R_i(\pi_i', \pi_{-i})$$

## 14.4 Nash Equilibrium vs Maximin

There was an important connection in perfect information between the solution concepts of Nash equilibria and maximin. While these concepts are at a first glance very different, in the case of zero-sum two player games, they collapsed (Theorem 5.1).

**Algorithm 4** Imperfect Information Game Best Response

1: **function** COMPUTEREACH($s \in S_i, i \in \mathcal{N}$ )
2:      **for** $h : \mathcal{S}_i(h)$ **do**
3:          **for** $h' = (haw) : \mathcal{H}$ **do**
4:                                          $\triangleright$ haw corresponds to some sas'
5:             h_reach_prob[$h'$] = h_reach_prob[$h$] $\cdot \mathcal{T}(h, a, h') \cdot \pi_{-i}(a)$
6:      **for** $sao' : \mathcal{S}_i$ **do**
7:          s_reach_prob[$sao'$] = $\sum_{h \in \mathcal{H}(sao')}$h_reach_prob[$h$]
8:          ComputeReach($sas'$, i)

9:                        $\triangleright$ Return best response value of the state $s$.
10: **function** BESTRESPONSEDFS($s \in \mathcal{S}_i, i \in \mathcal{N}$)
11:                          $\triangleright$ Terminal state, return terminal utility.
12:      **if** $\mathcal{A}_i(s) = \emptyset$ **then**
13:          **return** $R_i(s)$
14:        $\triangleright$ Action value is the reach-weighted sum of possible child states.
15:      **for** $a : \mathcal{A}(s)$ **do**
16:          sa_values[$sa$] = $\sum_{sao \in \mathcal{S}_i}$ BestResponseDFS($sao$) $\cdot$ s_reach_prob[$sao$]

17:
18:                          $\triangleright$ Choosing an action with the best value
19:      $\pi(s) = argmax_{sa}$ sa_values[sa]
20:      **return** $max_{sa}$ sa_values[sa]

21:
22: **function** BESTRESPONSE($s \in \mathcal{S}_i, i \in \mathcal{N}$ )
23:      s_reach_prob[$s$] = 1
24:      **for** $h : \mathcal{H}(s)$ **do**
25:          h_reach_prob[$h$] = 1
26:      COMPUTEREACH($s$, $i$)
27:      BESTRESPONSEDFS($s$, $i$)

This connection crucially relied on the two-player zero-sum property and imperfect information brings no change to these results. First, the minmax theorem (Theorem 2) still holds as the original referenced proof is for matrix games. Proofs of Theorem 4 and Theorem 5 then remain the same as there were no assumptions specific to perfect information.

## 14.5   Difference To Perfect Information

Imperfect information brought little to no change to the concepts of best response, maximin and Nash equilibrium. The same concepts, motivations and properties of the optimal policies hold in imperfect information. Optimal policies are still guaranteed to exist and provide the same appealing guarantees (e.g. it guarantees a positive expected reward against any opponent if we get to play both positions, Corollary 5.2) The reasons to follow optimal policies in imperfect information are thus as appealing as in perfect information. We just need more powerful algorithms to compute the policies. This is true for offline algorithms, but even more so for online search algorithms.

### 14.5.1   Deterministic or Stochastic Policies

One of the reasons why computing an optimal policy is more challenging in imperfect information is the need for careful randomization. In perfect information, an optimal policy can always be deterministic. But consider the imperfect information game rock-paper-scissors. It is easy to verify that any deterministic policy receives a worst-case utility of $-1$. On the other hand, the best response value of the uniform policy is 0. Deterministic policies thus might not be sufficient to play optimally in imperfect information settings (Corollary 17.1).

**Corollary 17.1.** *An optimal policy in imperfect information games might have to be stochastic.*

**Level of Uncertainty and the Support Size**

There is a rather interesting connection between the level of uncertainty and the support size (Definition 2). In perfect information games, there is only a single state of the opponent consistent with the player's state. In other words, there is no uncertainty about the opponent, the opponent can only be in one state and thus there's always an optimal strategy with support size one.

In imperfect information, there is uncertainty about the state of the opponent; there are multiple states the opponent can be in given a player's state. It turns out there is direct connection between the level of uncertainty and the number of actions required under an optimal policy (Schmid et al., 2014).

## 14.6   Multiplayer Games

While Section 4.6 already discussed that leaving the comfort of two-player zero sum settings greatly complicates matters, Section 4.6.3 deferred an example of the equilibrium selection problem. This is because there is a particularly nice

matrix game that illustrates this issue — the Game of Chicken (Rapoport and Chammah, 1966). It has been inspired by another famous game —- Prisoner's dilemma (Rapoport et al., 1965), and is now often re-branded as the "crossroad game" where the both players are allowed to either (S)top at the crossroad or (G)o. If both players stop, they both receive zero reward. If one of the players stop and the other goes, the cautious player receive a zero reward while the daring player receives positive reward. But if both players go at the same time, they crash and both receive a large negative reward. See Table 14.1 for the corresponding matrix game.

|          | (S)top | (G)o |          | (S)top | (G)o |
|----------|--------|------|----------|--------|------|
| (S)top   | 0      | 0    | (S)top   | 0      | 1    |
| (G)o     | 1      | -100 | (G)o     | 0      | -100 |

Table 14.1: Chicken game. Left: Player 1 utility. Right: Player 2 utility.

There are three Nash equilibria in this game:

1. Row player S, column player G

2. Row player G, column player S

3. Both players share the same strategy: $p(S) = \frac{100}{101}, p(G) = \frac{1}{101}$

The equilibrium selection problem is that the column player can select the first equilibrium, while the row player selects the second one. In this case, both players select the (G)o action and crash, which violates reward guarantee, and is not a Nash equilibrium.

# 15. Sub-games

This chapter generalizes sub-games in imperfect information settings. Just as in perfect information, a sub-game is going to be a well-defined game, a sub-problem of the original one based on the notion of sub-tree.

But while a sub-game in perfect information was simply defined by a sub-tree rooted in the current state (Section 5.1), things are more complicated in imperfect information. As the players have different views of the game, the possible states of the players differs and so do their respective (infostate) trees. There are two conceptual generalizations for sub-games in imperfect information games.

- Set of states we need to reason about.

- Distribution over the possible initial states.

## 15.1 Set of States to Reason About

Given a player's state, what is the set of states we need to reason about? We must consider all the states the other player can potentially be in. In imperfect information games, there might be multiple such consistent states consistent with the current observation. And to reason about any of these states, the opponent has to now consider all the possible states the first player could be in. This recursive reasoning continues until we find the set of all relevant states for the current situation.

### 15.1.1 Consistent States

To formalize the set of states we need to reason about, we first define consistent states — the set of (opponent's) states consistent with the current state (Definition 21). The corresponding operation $\mathcal{C} : \mathcal{S} \to 2^{\mathcal{S}}$ is then referred to as the consistency operation.

**Definition 21** (Consistent States). *Given a state $s \in \mathcal{S}_i$, the set of (opponent's) consistent states $\mathcal{C}(s)$ is:*

$$\mathcal{C}(s) = \{s' \in \mathcal{S}_{-i} \,|\, \mathcal{H}(s) \cup \mathcal{H}(s') \neq \emptyset\} \tag{15.1}$$

### 15.1.2 Common Information Set

The consistency operation represents a single step of the recursive process. We need to repeat the operation until we find a closure — set of states closed under the consistency operation. This particular closure is referred to as the common information set (Definition 22). The common information set represents the minimal set of relevant states.

**Definition 22** (Common Information States). *Given a state $s$, the common information set is the closure under the consistency operation: $cl_{\mathcal{C}}(s)$.*

### 15.1.3 Common vs Public Information

While common information defines the minimal set of relevant states, there is an easier way to find a set closed under the consistency operation without the need for the recursive construction process (closure). We just need to consider all the states sharing public information.

As the common information set is a subset of the public state set: $cl_{\mathcal{C}}(s) \subseteq \mathcal{S}_{pub}(s)$, the set of states sharing a public information is closed under the consistency operation $\mathcal{C}$. The only downside of the public state is that it can potentially be strictly larger. In many games though (e.g. in all example games), $cl_{\mathcal{C}}(s)$ and $\mathcal{S}_{pub}(s)$ are exactly the same and thus the notion of public states is preferable due to its simplicity.

### 15.1.4 Future States

Introduced sets described the necessary states relevant for the current situation. In a sub-game, we must reason about the current states as well as all the reachable future states. Given state $s$, we need to reason about $\mathcal{S}_{pub}(s)$ and all the future states of this set $\{f | s' \sqsubset f, s' \in \mathcal{S}_{pub}(s)\}$. But do we need to repeat the construction of consistent states for these future states? Given a future state $f$, are there any more relevant states that are not already included? Theorem 18 guarantees that this is not the case.

**Theorem 18.** *For any future states $f$ of a state $s$:*

$$s \sqsubset f, f' \in \mathcal{S}_{pub}(f) \implies \exists s' \in \mathcal{S}_{pub}(s) : s' \sqsubset f'$$

*Proof.* Follows from the perfect recall property — there is a unique trajectory from a root of the game to each infostate (infostates form a directed tree). □

All the states we need to reason about are thus the states in a public state and all its children — a public sub-tree.

### 15.1.5 Public Subtree

We have now finished the first essential piece of the subgame generalization. In perfect information, the set of relevant states was simply the sub-tree rooted in the current state. In imperfect information, we need to use the public state sub-tree.

Note that in the case where common information would be preferable to public information (the corresponding set being substantially smaller), all the results apply to the sub-games formed by common information.

### 15.1.6 Walk-through Example

To get some more intuition about the notions of consistency, common and public information sets as well as the sub-trees, we will illustrate these concepts on a simple poker game with asymmetrical dealing of cards presented in the Figure 15.1.

Figure 15.1: Simple poker game with a three card deck {Q, K A}, but the possible deals for the first and second player are {(A, A), (K, A), (K, K), (Q, Q)}. This asymmetrical deal distribution helps to illustrate the details of common and public information sets.

## Common Information States

We start with the concept of consistent states and the corresponding common information set, and we will consider two different initial states: the first player being dealt Queen or Ace.

**First Player Dealt Queen**  The state we start our recursive construction is $s = [d|Q|]$ - see also Figure 15.2.

1. Common information states start with $cl_{\mathcal{C}}^0(s) = \{[d|Q|]\}$.

2. $\mathcal{C}([d|Q|]) = \{[d||Q]\}$, expanding the set $cl_{\mathcal{C}}^1(s) = \{[d|Q|], [d||Q]\}$.

3. $\mathcal{C}([d||Q]) = \{[d|Q|]\}$, ending our recurrent reasoning.

4. $cl_{\mathcal{C}}(s) = \{[d||Q], [d|Q|]\}$

**First Player Dealt Ace**  The state we start our recursive construction is $s = [d|A|]$ - see also Figure 15.3.

1. Common information states start with $cl_{\mathcal{C}}^0(s) = \{[d|A|]\}$.

2. $\mathcal{C}([d|A|]) = \{[d||A]\}$, expanding the set $cl_{\mathcal{C}}^1(s) = \{[d|A|], [d||A]\}$.

3. $\mathcal{C}([d||A]) = \{[d|A|], [d|K|]\}$, and $cl_{\mathcal{C}}^2(s) = \{[d|A|], [d||A], [d|K|]\}$.

4. $\mathcal{C}([d|K|]) = \{[d||K]\}$, and $cl_{\mathcal{C}}^3(s) = \{[d|A|], [d||A], [d|K|], [d||K]\}$.

5. $\mathcal{C}([d|K|]) = \{[d||K]\}$, ending our recurrent reasoning.

6. $cl_{\mathcal{C}}(s) = \{[d|A|], [d||A], [d|K|], [d||K]\}$.

74

(a) Infoset $[d|Q]$ of the first player.



(b) Infoset $[d|Q]$ of the second player.

Figure 15.2: (a) When reasoning about the player's state $[d|Q]$, we need to reason about all the possible current and future states of both players - and continue this reasoning process until we get a closure. (b) First level of such recursive reasoning.

(a) $cl_{\mathcal{C}}^0(s) = \{[d|A|]\}$



(b) $\mathcal{C}([d|A|]) = \{[d||A|])\}$, $cl_{\mathcal{C}}^1(s) = \{[d|A|]), [d||A]\}$



(c) $\mathcal{C}([d||A]) = \{[d|K|])\}$, $cl_{\mathcal{C}}^2(s) = \{[d|A|]), [d||A], [d|K|]\}$



(d) $\mathcal{C}([d|K|]) = \{[d||K|])\}$, $cl_{\mathcal{C}}^2(s) = \{[d|A|]), [d||A], [d|K|], [d||K]\}$

Figure 15.3: $s = [d|A|]$

**Public States**

There is no need for the recurrent construction when public states are used, we simply use the corresponding public state $P(s)$ of the current state. In this particular examples, both initial states we considered (player dealt Queen or Ace) are within the same public state (Figure 15.4) and the public state is larger than the common information sets. This is because of the particular choice of the asymmetrical deal. In typical poker variants, the sets based on common and public information are identical.



Figure 15.4: Two initial states $s_1 = [d|A|]$, $s = [d|A|]$, $s_2 = [d|Q|]$ share the same public states. In this example, public state is larger than the common information sets.

## 15.2 Distribution

The set of states formed by a public sub-tree identifies the states we can reason about in isolation. In perfect information, a strategy for all the states in a sub-tree was sufficient for value computation of these states. That is not the case in imperfect information, where the transition dynamics and consequently the utilities depend not just on the strategy in the sub-tree, but also on the strategy above. We have already seen this property in the best response calculation (Section 14.1.1 and Algorithm 4). Concretely, the transition dynamics in a sub-tree are fully identified by the distribution over the possible initial histories and a strategy for all the infostates of that sub-tree.

### 15.2.1 Factorized Distribution

Computing values in a sub-tree requires a strategy for all such states as well as a distribution over the initial histories. But not any distribution over those histories can result from an agent's play as they have control only over actions in their respective infostate (rather than for histories).

Given an initial public state, the distribution over the histories is a function of i) fixed chance factor and ii) players' strategies above the public state. The chance factor is fixed and can be different for any history. Agent strategies on

the other hand are not fixed and are the same for all histories grouped in a single infostate.

For a history $h$, the reach probability is a product of the chance contribution $P_c(h) = \Pi_{h'aw \sqsubseteq h} \mathcal{T}(h', a, w)$ and players' reach probabilities: $P^\pi(h) = \Pi_{i \in \mathcal{N}} P_i^{\pi_i}(h)$. A player's reach probability $P_i^{\pi_i}(h)$ is then a product of the strategies in their infostate sequence leading to the that history: $P_i^\pi(h) = \Pi_{h'aw \sqsubseteq h} \pi_i(s_i(h), a)$. As this quantity is equal for all the histories in the infostate $s_i(h)$, we have $P_i^\pi(s_i(h)) = P_i^\pi(h)$.

We can thus compute $P^\pi(h) = P_c(h) \Pi_{i \in \mathcal{N}} P_i^{\pi_i}(s_i(h))$, factoring this term to each player's infostate reach probabilities. In order to compute a distribution over the histories in a public state, we only need to know reach probabilities for all the players' infostates in that public state. The importance of this observation is that the number of infostates is often substantially smaller.

Consider Texas Hold'em poker as an illustrative example. The number of player's infostates in a public state is $\binom{52}{2}$ (player is dealt two private cards), while the number of individual histories is $\binom{52}{4}$ (two private cards are dealt to each player). To represent the distribution over all these histories, we only require reach probabilities for the infostates of both players: $2\binom{52}{2}$, which is more than 100 times smaller from $\binom{52}{4}$.

## 15.3   Public Sub-game

We now have all we need for the most critical concept for search — generalization of sub-games in imperfect information. While in perfect information, a sub-game was identified simply by a sub-tree, imperfect information sub-games are identified by:

1. Public state $s_{pub} \in \mathcal{S}_{pub}$

2. Per-player distribution over their info-states $\Delta(\mathcal{S}_i(s_{pub}))$.

As we will see, this sub-game allows for all the same decomposition concepts as in perfect information. Namely, it allows for:

- Solving for optimal policies in sub-games.

- Defining value functions for sub-games.

- Decomposition.

- Construction of provably safe search methods.

### 15.3.1   Public Sub-Game is a Game

Simple construction reveals that the public sub-game forms a well-defined imperfect information game. The ideas is that we construct a game as if it started in the public sub-tree with the appropriate distribution over the initial states. We simply construct a game rooted in a chance node that "deals out" the initial histories/states , with probability of each history being the (normalized) product of each player's contribution (including chance) — Definition 23.

| | Input | Output |
|---|---|---|
| Perfect Information | $s \in \mathcal{S}$ | $V(s)$ |
| Imperfect Information | $s_{pub} \in \mathcal{S}_{pub}, \Delta(\mathcal{S}_1(s_{pub})), \Delta(\mathcal{S}_2(s_{pub}))$ | $V(s) \, \forall s \in \mathcal{S}_i(s_{pub})$ |

Table 15.1: Value function in perfect and imperfect information settings.

**Definition 23** (Public Sub-game)**.** *Public sub-game defined by a tuple ($s_{pub} \in \mathcal{S}_{pub}, \Delta(\mathcal{S}_1(s_{pub})), \Delta(\mathcal{S}_2(s_{pub}))$) is a game rooted in a chance state $w_0$ with actions dealing out initial histories: $\mathcal{A}_0(w_0) = \mathcal{H}(s_{pub})$, and the strategy (distribution): $\mathcal{T}(w_0, h) \propto \Delta(\mathcal{S}_1(h)) \, \Delta(\mathcal{S}_2(h)) \, P_c(h)$*

## 15.4 Value Functions

Analogously to value functions for perfect information games, value functions map sub-games to values. And as a public sub-game is a game, we also use the value under an optimal policy.

An important generalisation is that in perfect information, the value of a sub-game was a single scalar. In imperfect information, it is a vector of values, a value for each initial infostate (Table 15.1). As nothing limits the imperfect information case to contain a single state of both players, we also get Observation 18.1.

**Observation 18.1.** *Perfect information value functions are a degenerate case of imperfect information value function with a single initial state*

### 15.4.1 Example

Table 15.2 illustrates the value function in rock-paper-scissors. Consider the sub-game irooted in the public state just after the first player acted and the second player is to act. That public state contains a single infostate of the second player and three infostates of the first player (after the rock, paper or scissors action). The sub-game and the value function thus requires a distribution over the three possible initial states of the first player, and the probability of the second player being in their single state is 1. The output then consists of a single value for each of the players' states (three values for the first player and one value for the second player).

### 15.4.2 Unique Game Value, Multiple State Values

While the game value is unique and thus the same under any optimal policy, individual state values can differ. Different optimal policies can produce different infostate values. In the rock-paper-scissors example, any policy of the second player is optimal under the uniform initial distribution of the first player. This in turn corresponds to different state values of the first player (Table 15.2).

But why do we care about these individual state values at all? Why is not the game value sufficient? Its importance will become clear once these value function are used within search, where we will need the individual state values.

| Input | | Output | | Optimal Policy |
|---|---|---|---|---|
| $\Delta(\mathcal{S}_1(s_{pub}))$ | $\Delta(\mathcal{S}_2(s_{pub}))$ | $V_1$ | $V_2$ | |
| $(0.2, 0.2, 0.6)$ | $(1)$ | $(0, 1, -1)$ | $(0.4)$ | $\pi_2 = (1, 0, 0)$ |
| $(0.4, 0.3, 0.3)$ | $(1)$ | $(-1, 0, 1)$ | $(0.1)$ | $\pi_2 = (0, 1, 0)$ |
| $(\nicefrac{1}{3}, \nicefrac{1}{3}, \nicefrac{1}{3})$ | $(1)$ | $\pi_2 \begin{pmatrix} 0 & 1 & -1 \\ -1 & 0 & 1 \\ 1 & 0 & -1 \end{pmatrix}$ | $(0)$ | any policy $\pi_2$ |

Table 15.2: Value function in rock-paper-scissors for the sub-game rooted in the second public state. Note that as the sub-game is zero-sum, the game values are in balance $\Delta(\mathcal{S}_1(s_{pub}))V_1^\top + \Delta(\mathcal{S}_2(s_{pub}))V_2^\top = 0$. Value for the first player in the sub-game is maximized when the initial beliefs are uniform $(\nicefrac{1}{3}, \nicefrac{1}{3}, \nicefrac{1}{3})$.

| Input | | Output | | Optimal Policy |
|---|---|---|---|---|
| $\Delta(\mathcal{S}_1(s_{pub}))$ | $\Delta(\mathcal{S}_2(s_{pub}))$ | $V_1$ | $V_2$ | |
| $(\nicefrac{1}{3} + 2\epsilon, \nicefrac{1}{3} - \epsilon, \nicefrac{1}{3} + \epsilon)$ | $(1)$ | $(0, 1, -1)$ | $(\epsilon)$ | $\pi_2 = (0, 1, 0)$ |
| $(\nicefrac{1}{3} - 2\epsilon, \nicefrac{1}{3} + 4\epsilon, \nicefrac{1}{3} - 2\epsilon)$ | $(1)$ | $(-1, 0, 1)$ | $(2\epsilon)$ | $\pi_2 = (0, 0, 1)$ |
| $(\nicefrac{1}{3} - 3\epsilon, \nicefrac{1}{3} - 3\epsilon, \nicefrac{1}{3} + 6\epsilon)$ | $(1)$ | $(1, -1, 0)$ | $(3\epsilon)$ | $\pi_2 = (1, 0, 0)$ |

Table 15.3: Small perturbations in the input distribution can lead to large changes in the underlying optimal policy and thus output state values, but to only small changes in the game value.

### 15.4.3 Structure

As the value function in imperfect information is more complex due to the beliefs, there are some important questions related to the structure of the function. Do small perturbations of the input distribution lead to drastic changes in the output? Is there some structure? The structure of the function becomes even more relevant once we get to function approximation and learning (e.g. with neural networks), as learning/generalization with no structure is not possible (Wolpert, 1996).

The game value itself ($\Delta S_i(s_{pub})V_i^\top$) has a particularly nice, piece-wise linear structure. See Section 16.1.1 and Figure 16.1 for some graphical intuition, or e.g. Čermák et al. (2017); Seitz et al. (2019); Kovařík et al. (2020); Brown et al. (2020) for more analysis and properties. But how about the individual state values $V_i$? In general, small input perturbations of the input distribution can lead to large changes of the output state values (Table 15.3). In other words, small perturbation in the definition of the game can lead to drastic changes in how the game is played optimally, but only to small change in the game value.

Fortunately, this sensitivity happens only at the critical points of the otherwise piece-wise linear value function (Figure 16.1). Furthermore, it is possible to bound the state value changes when the underlying policy is not optimal, but rather $\epsilon$-optimal (Theorem 19).

**Theorem 19.** *Let $(V_1^a, V_2^a)$ be the state values of a sub-game $G^a = (s_{pub}, d_1^a \in \Delta(\mathcal{S}_1(s_{pub})), d_2^a \in \Delta(\mathcal{S}_2(s_{pub})))$ under an optimal policy $\pi^a$. Le $d_1^b, d_2^b$ be $\epsilon$-perturbed*

*distributions ($|d_i^a - d_i^b|_\infty < \epsilon$). There is an $\epsilon_2$-optimal policy $\pi^b$ for the perturbed sub-game $G^b = (s_{pub}, d_1^b, d_2^b)$ producing $\epsilon_2$-perturbed state values $(V_1^b, V_2^b)$ where $\epsilon_2 < \epsilon\Delta_R$.*

*Proof.* Following the original policy $\pi^a$ in the perturbed game is an $\epsilon_2$-optimal policy producing $\epsilon_2$ perturbed values, as the change in distribution over the terminal states is bounded by $\epsilon$. □

# 16. Offline Solving I

We will now describe offline solving methods for imperfect information games. Because we include multiple distinct approaches, we split the methods into three separate chapters.

Chapter 16 deals with optimization methods and self-play methods. Section 16.1 starts with matrix games, and we analyse the structure of the worst case value as a function of the agent's policy (Section 16.1.1). This is then used in construction of a linear optimization problem that produces an optimal policy (Section 16.2) and allows us to prove the minimax theorem. Section 16.3 generalizes the construction of the linear optimization problem to the sequence games and formalizes the sequence form. Section 16.4 discusses the application of (single agent) reinforcement learning methods, showing that naive application of such methods in self-play fails to converge even in simple imperfect information games. Section 16.5 touches on some sound multi-agent reinforcement learning methods and Section 16.6 introduces fictitious play. Double oracle and its related methods are then mentioned in Section 16.7.

Chapter 17 moves to the most important offline solving methods of this thesis — methods based on the regret minimization framework that are central for the final search techniques. Section 17.1 presents the online convex learning concept of regret minimization, and Section 17.2 includes some regret minimization algorithms. Section 17.3 then shows the striking connection between regret minimization and optimal policies, where we prove the Folk Theorem (22) that shows that if both players use a regret minimizer in self-play settings, the average policies converge to an optimal policy profile. Section 17.5 describes how to minimize regret in sequential settings — counterfactual regret minimization (CFR). CFR decomposes the full regret into a sum over individual regrets at each information state, allowing one to independently minimize these partial regrets. Section 17.6 then mentions some modern members of the CFR algorithm family, such as CFR+ and DCFR. Section 17.7 describes Monte Carlo variants of CFR (MC-CFR) that sample trajectories rather than traversing the entire game tree at each iteration. As the sampling introduces noise and slows down convergence, Section 17.7 describes VR-MCCFR, modern variance reduction methods that drastically speed up the convergence (a contribution of this thesis).

Chapter 18 deals with large games where we cannot directly use tabular methods. Section 18.1 discusses abstraction approaches that were for a long time state of the art for large imperfect information games (until the dawn of search). Section 18.2 then briefly mentions some other non-tabular methods that can tackle large games.

## 16.1  Direct Optimization for Matrix Games

An optimal policy maximizes its best response value $\max_{\pi_i} BRV(\pi_i)$. For matrix games, the $BRV(\pi_i)$ function has a particularly simple form (16.1).

$$BRV(\pi_i) = \min_{\pi_{-i}} \pi_i A \pi_{-i} \tag{16.1}$$

Since 16.1 is a point-wise minimum of a affine function, the function is concave (Boyd and Vandenberghe, 2004).

### 16.1.1 Graphical Interpretation

To build more intuition about the concave best response value function, we will graphically analyze the function on a small matrix game. Table 16.1 shows a game with 2 actions for the row player, and 3 actions for the column player.

|   | X | Y | Z |
|---|---|---|---|
| A | 8 | 3 | 2 |
| B | 2 | 6 | 8 |

Table 16.1: Simple game for the analysis of the best response value function.

Visualising the game in a particular way (Figure 16.1a) provides nice insight into the structure of the best response value function. We can see that the function is concave and piece-wise linear — key observation for algorithms that directly maximize the function. This graphical reasoning can be traced to Gale et al. (1951).



(a) Visualization of Game 16.1.

(b) Worst case function we try to maximize.

(c) Strategy that maximizes the worst case function.

Figure 16.1: a) We can see the three lines corresponding to the value of the row play when the opponent responds with one of the three available actions. The x-axis corresponds to the strategy of the row player. As an example, the depicted red point is the return for the row player when they play $\pi_1 = (0.2A, 0.8B)$ and the opponent responds with action Z. b) Best response value function is then simply minimum of the lines (worst case best response). c) Strategy maximizing the worst case function.

## 16.2 Matrix Games Linear Programming

The piece-wise linear structure of 16.1 can be leveraged to create a closed-form linear program that optimizes the function in question. We will construct a linear optimization for both the row and the column players, and then use the strong duality to prove the minmax theorem.

## 16.2.1 Row Player

First, notice that maximin solution for the row player is:

$$\max_{x \in \Pi_i} \min_{y \in \Pi_{-i}} x^\top A y \tag{16.2}$$

Since (16.2) is bi-linear, we need to decouple the $x$ and $y$ by introducing new variable. Given a strategy $x$ for the row player, the best-responding opponent simply chooses the column with the smallest utility.

$$\min_{y \in \Pi_{-i}} x^\top A y \tag{16.3}$$

This can be re-formulated as:

$$\max_{u \in \mathbb{R}} u \\ x^\top A \geq u \tag{16.4}$$

We now just use this formulation in (16.2), resulting in the linear program (16.5) and Theorem 20.

$$\max_{u \in \mathbb{R}, x \in \Pi_i} u \\ x^\top A \geq u \tag{16.5}$$

**Theorem 20.** *(Nisan et al., 2007, Theorem 1.11) Linear program 16.5 produces a maximin strategy for the row player (x) and the corresponding maximin value (u).*

## 16.2.2 Column Player

Analogical reasoning for the column player starts with $\min_y \max_x x^\top A y$ and results in (16.6) and Theorem 20.

$$\min_{v \in \mathbb{R}, y \in \Pi_{-i}} v \\ A y \leq v \tag{16.6}$$

**Theorem 21.** *Linear program 16.6 produces a maximin strategy for the column player (y) and the corresponding maximin value (−v).*

## 16.2.3 Minimax Theorem

We are now ready to prove the crucial minimax theorem.

**Theorem 2** (Minimax Theorem)**.**

$$\max_{\pi_i} \min_{\pi_{-i}} R_i(\pi_i, \pi_{-i}) = \min_{\pi_{-i}} \max_{\pi_i} R_i(\pi_i, \pi_i) \tag{4.10}$$

The key observation is that the linear programs for the maximin strategies of the row player 16.5 and the column player 16.6 form a pair of dual linear programs. And since both are feasible, strong duality holds and thus the optimal values of these programs are equal $v = -u$ (Bertsimas and Tsitsiklis, 1997).

It is not a coincidence that we have used linear programming and duality to find optimal policies in two-person zero-sum games. There is a one-to-one connection between linear programming and solving optimal policies, as it is possible to convert one problem to another (Boyd and Vandenberghe, 2004).

## 16.3   Sequential Games Linear Programming

This section generalizes the linear programming formulation of optimal policies from matrix games to sequential decision making. Even though it is possible to convert a sequential game to matrix form, this is not practical as the resulting matrix game can be exponentially large compared to the sequential representation (Section 13.2.3).

### 16.3.1   Sequence Form

While a behavioral strategy defines how to act in individual information states, this representation is not directly useful for linear/convex optimization due to its non-convexity (probability of reaching a state $s$ is a product of the behavioral probabilities on its infostate path $\Pi_{s'as'' \sqsubseteq s} \pi(s', a)$) Fortunately, there is a convex representation of a strategy — the sequence form. The idea is to represent a probability of a player reaching a state as the product of the sequence of a player's actions leading to that state. One then must make sure that these reach probabilities are well formed, that the probability mass of reaching all children state sum to their parent.

In this section, we present a sequence form representation in the FOSG formalism by adapting the notation similar to that for extensive form games in Nisan et al. (2007).

**Realization Plan**

We use $\mathcal{B}_i$ to denote the set of all state-action sequences of a player $i$: $\mathcal{B}_i = \{sa | s \in \mathcal{S}_i, a \in \mathcal{A}(s)\}$, and $b_i(ha) = s_i(h)a$. The sequence form realization probability of a state-action sequence $sa \in \mathcal{B}_i$ becomes $\prod_{s'a' \sqsubseteq sa} \pi_i(s', a')$. For game with perfect recall, Kuhn's theorem (Kuhn, 1953; Aumann, 1961) states that realization plan is equivalent to behavioral strategy.

We use $x$ and $y$ to denote the vectors of these quantities for all such sequences of player 1 and 2 respectively. Recall that the terminal utility (return) of reaching $z \in \mathcal{Z}$ is simply the accumulated reward on that trajectory:

$$u_i(z) = R_i(z) = \sum_{haw \sqsubseteq z} R_i(h, a, w) \qquad (16.7)$$

We are now ready to define the sequence form payoff matrix $A \in \mathbb{R}^{|\mathcal{B}_1| \times |\mathcal{B}_2|}$, with $A_{\sigma\tau}$ for indexes $\sigma \in \mathcal{B}_1, \tau \in \mathcal{B}_2$ defined as the chance-reach-weighted sum over the corresponding terminal states:

$$A_{\sigma\tau} = \sum_{(haw) \in Z \,:\, b_1(ha) = \sigma, b_2(ha) = \tau} P_{\mathcal{T}}(z) u_1(z) \tag{16.8}$$

With this notation, the expected value given realization plans of both players becomes simply:

$$x^\top A y \tag{16.9}$$

Not just any positive real valued vectors form a realization plan though. A realization plan fulfills two important properties (ensuring that the probability mass of children sequences sum to their parent):

$$x_\emptyset = 1 \tag{16.10}$$

$$\sum_{a' \in \mathcal{A}(sao)} x_{saoa'} = x_{sa} \tag{16.11}$$

These conditions can be compactly formulated using zero-one matrices $E, F$ and zero-one vectors $e, f$ (Nisan et al., 2007):

$$Ex = e, \; x \geq \mathbf{0} \tag{16.12}$$
$$Fy = f, \; y \geq \mathbf{0} \tag{16.13}$$

### 16.3.2 Optimal Policies

We will now use the sequence form to construct a linear optimization problem for computing an optimal policy. Given a fixed realization plan $y$ of the second player, we can compute the best response for the first player as:

$$\max x^\top A y$$
$$Ex = e, \; x \geq \mathbf{0} \tag{16.14}$$

This can be turned into the following LP that solves for Nash equilibrium (for more details see Nisan et al. (2007)).

$$\min_{u,v} e^\top u$$
$$Fy = f, E^\top u - Ay \geq \mathbf{0}, y \geq \mathbf{0} \tag{16.15}$$

## 16.4 Self-play via Independent Reinforcement Learning

We now consider another approach for producing optimal offline policies — methods rooted in the self-play paradigm. The first such method is self-play via independent reinforcement learning, where the agents best-respond to each other at each time step. In perfect information, this sequence[1] converged to an optimal policy (Section 6.2). The same algorithm fails to produce optimal policies in imperfect information games. Table 16.2 presents a simple matrix game with optimal policies $\pi_1 = (\frac{2}{3}, \frac{1}{3})$ and $\pi_2 = (\frac{2}{3}, \frac{1}{3})$ for the row and column player respectively. A best-responding sequence in this game $p_1 = \{A, B, A, B \ldots\}$, $p_2 = \{B, A, B, A \ldots\}$ then does not converge as the strategy at each time $p_i^t$ is highly exploitable. Furthermore, note that the average strategies then are $\pi_1 = (\frac{1}{2}, \frac{1}{2})$ and $\pi_2 = (\frac{1}{2}, \frac{1}{2})$ and thus also exploitable.

|   | A | B |
|---|---|---|
| A | 1 | -2 |
| B | -2 | 4 |

Table 16.2: Matrix game where best-responding sequence does not converge.

While the same argument could be made using the rock-paper-scissors, the example would be far less illustrative. The problem is that the relative frequencies of the individual actions in the best responding sequence $p_1 = \{R, P, S, R, P, S, \ldots\}$, $p_2 = \{R, P, S, R, P, S, \ldots\}$ match the optimal policy. One could then conclude that average strategy converges, but that is not generally the case.

### 16.4.1 Deterministic Optimal Policies

As optimal policies under imperfect information might have to be stochastic (Corollary 17.1), the failure of the selfplay method to converge could be attributed to the fact that it produces a deterministic best response policy at each time step. But a best-responding sequence can fail to converge even in a game with an optimal deterministic policy. Table 16.3 presents a modified rock-paper-scissors game with added action (W)ater. Water draws against all the actions, and both players playing water is an optimal policy in this game (there are other optimal policies). Sequence $p_1 = \{R, P, S, R, P, S, \ldots\}$, $p_2 = \{R, P, S, R, P, S, \ldots\}$ is then a best-responding sequence failing to converge in this game.

## 16.5 Multi-agent Reinforcement Learning

The message of the convergence failure should be that we just cannot naively take a single-agent reinforcement learning algorithm and run it independently in multi-agent settings through self-play. Single-agent reinforcement learning methods lack theoretical guarantees in these settings, often producing highly exploitable policies. While we have simplified the analysis by allowing the algorithm

---

[1]Sub-game perfect best response.

|     | R   | P   | S   | W   |
| --- | --- | --- | --- | --- |
| R   | 1   | 0   | -1  | 0   |
| P   | -1  | 1   | 0   | 0   |
| S   | 0   | -1  | 1   | 0   |
| W   | 0   | 0   | 0   | 0   |

Table 16.3: Rock-paper-scissors-water: game where deterministic optimal policy exists (both player playing water), but best-responding selfplay sequence can fail to converge.

to fully converge to best response at each self-play step, this approach fails to produce optimal policies even if the agents improve their policies less aggressively (e.g. using few steps of policy gradient or Q-learning). For more insights and analysis see e.g. Bailey and Piliouras (2018). Despite the lack of theoretical guarantees, this approach is relatively popular and it has recently been applied to e.g. Dota 2 (Berner et al., 2019). As such approach fails to produce optimal policies for rock-paper-scissors, it is a questionable choice for large two-player imperfect information game.

Rather than naively applying single-agent algorithms, one can properly design the algorithms for multi-agent settings. There are many modern methods that take the complex multi-agent dynamics into account, providing strong convergence guarantees (Muller et al., 2019; Vinyals et al., 2019a; Perolat et al., 2020) None of these methods allow for search nor provide any building blocks for the next chapters and thus we do not dive any deeper into details.

## 16.6   Fictitious Play

Simple modification to the best-responding sequence leads to the fictitious play algorithm. Fictitious play best-responds to the average strategy of the opponent rather than the last one.

Let $\overline{\pi_i^t}$ denote the average strategy of the player $i$ up to time $t$. The fictitious self-play sequence is then $\pi_i^t \in \mathbb{BR}(\overline{\pi_{-i}^{t'-1}})$. Best-responding to the average policy has important theoretical consequences, as the average strategy profile then converges to optimal. This algorithm dates all the way back to (Brown, 1951), with a convergence proof followed shortly after (Robinson, 1951).

Fictitious play is very simple to implement and often serves as a building block for other algorithms. For example in large games, one can use single-agent reinforcement learning methods to approximate the individual best responses (Heinrich et al., 2015). Weakened fictitious play then replaces the best-response to the average strategy with a strategy that is only an increasingly better one (Van der Genugten, 2000; Leslie and Collins, 2006) The problem of this algorithm is its convergence speed, where the number of iterations required can grow exponentially (Harris, 1998).

## 16.7 Double Oracle Methods

Rather than using best response to produce the self-play sequence, it can be used to iteratively build and expand the matrix game to be solved. Double oracle methods (McMahan et al., 2003) iteratively expand the subset of actions to be considered by both the row and column player at time $t$ : $\mathcal{A}_1^t \subseteq \mathcal{A}_1$, $\mathcal{A}_2^t \subseteq \mathcal{A}_2$. This restricted matrix game $\mathcal{A}_1^t, \mathcal{A}_2^t$ is then solved e.g. via linear optimization. Next, both players compute a best response against the resulting optimal policy by considering all the actions in the original unrestricted game $br_i^t \in \mathcal{A}_i$. The restricted game is then expanded $\mathcal{A}_i^{t+1} = \mathcal{A}_i^t \cup br_i^t$ and the process repeats. Double oracle methods terminate when the restricted game does not expand any more ($\mathcal{A}_i^{t+1} = \mathcal{A}_i^t$) and the key result is that the optimal policy for the restricted game is then also optimal for the original game. This approach often terminates when the restricted game is substantially smaller than the original game. But in the worst case, the restricted game fully recovers the original game (even when the support size is smaller).

Policy-Spaced Response Oracles (PSRO) then extends this idea for large games where the best response calculation is intractable, and uses reinforcement learning methods to approximate the best response (Lanctot et al., 2017). The restricted game consists of individual policies represented by neural networks and is referred to as the meta-game. There are now multiple parallel and efficient variants of this approach (McAleer et al., 2020). Furthermore, AlphaStar (an AI for the game StarCraft II) builds on similar methods and the resulting meta-game with 888 rows and columns has been released (Vinyals et al., 2019a,b).

Finally, double oracle methods can be efficiently generalized to sequential settings (Bosansky et al., 2014). Combination with PSRO results in efficient PSRO methods in the sequential representations (McAleer et al., 2021).

# 17. Offline Solving II - Regret Minimization

Currently, the most successful algorithms for solving imperfect information games are based on the regret minimization framework. Regret minimization is a general, online convex learning concept where an agent repeatedly makes decision against an unknown environment (Zinkevich, 2003). Regret then measures the difference between the accumulated reward and the reward that a best time-independent action would receive in hindsight (e.g. always playing rock). And while the environment can be adversarial as the reward vector at each time step can depend on the selected action, it is still possible to provide strong guarantees with regard to hindsight performance of regret. An algorithm is said to be Hannan consistent if the regret grows sub-linearly — the average regret converging to zero.

Regret minimization has an elegant and important connection to games when used in self-play. In self-play, the reward vector can be computed using the opponent's policy and both players then select a next strategy using a regret minimizer. It is then possible to show that the average strategy of the players is $\epsilon$-optimal for $\epsilon$ no greater than the sum of players' average regrets. And as the average regret converges to zero, the average strategy converges to optimal.

Finally, minimizing regret in sequential settings can be decomposed into individual regret minimizers in the information states — counterfactual regret minimization. This is particularly important as it will be used as a key building block of decomposition and search methods.

## 17.1  The Setup

While there are multiple concepts of regret, the one used in this book is often referred to as external regret[1]. The agent repeatedly makes a decision against an unknown environment and observes a reward vector (reward for each action).

Formally, the agent's set of actions is $\mathcal{A}$ and at each time step $t$ the agent chooses a policy $\pi^t \in \Delta(\mathcal{A})$. The agent then receives a reward vector $x_t \in \mathcal{R}^{|\mathcal{A}|}$ and the agent's value is the weighted sum $v^t = \sum_{a \in A} \pi^t(a) \, x^t(a) = \pi^t \, x^{t\mathsf{T}}$. The cumulative reward up to time $T$ is simply $X_\pi^T = \sum_{t=1}^T \pi^t \, x^{t\mathsf{T}}$. External regret of action $a$, $R_a^T$ then contrasts this to a cumulative reward that would have been received if we rather followed action $a$ at each time step $R_a^T = \sum_{t=1}^T x_t(a) - X_\pi^T$. Finally, external regret is then $R^T = \max_{a \in \mathcal{A}} R_a^T$. In other words, how much we regret not taking the best action in retrospect.

## 17.2  Regret Minimization Algorithms

We now present some popular regret minimization algorithms. To simplify the notation and analysis, we assume that the action reward $x^t(a) \in [0, 1]$. No

---

[1]External regret with a comparison class of all the actions.

generality is lost as for $x^t(a) \in [a, b]$, the resulting bounds are simply scaled by a constant factor $|b - a|$

## 17.2.1 "Always Rock"

Selecting the same action $a_1$ at each time step leads to worst case $R^T = T$, as we can assign $x^t(a_1) = 0$, $x^t(a') = 1 \, \forall t$.

## 17.2.2 Simple Greedy

Instead of selecting the same action, we might prefer the action with the highest accumulated regret $\pi^t = \arg\max_{a \in \mathcal{A}} R_a^T$ (and selects the action with the smallest index if there are multiple maxima). It is easy to show that the regret can be bounded by (17.1).

$$R^T_{\pi_{greedy}} \leq \frac{T}{|\mathcal{A}|} \tag{17.1}$$

This is not much of an improvement as we need sub-linear regret for Hannan consistency. Turns out that in order to guarantee sub-linear regret, it is necessary to produce stochastic policies $\pi^t$ (Nisan et al., 2007).

## 17.2.3 Hedge

Arguably the most popular regret minimization algorithm — Hedge (Freund and Schapire, 1997) — is a parametric algorithm based on the weighted majority algorithm (Littlestone and Warmuth, 1994), selecting an action exponentially proportional to its regret (17.2).

$$\pi_t = \frac{e^{\beta R_{t-1}}}{\sum_{a \in A} e^{\beta R_{t-1}}} \tag{17.2}$$

Given the sequence length $T$ ahead of time and setting $\beta = \sqrt{\frac{2 \log(|\mathcal{A}|)}{T}}$, Hedge provides sub-linear upper bound on the regret (17.3).

$$R^T(\pi_{hedge}) \leq \Delta\sqrt{2T \log(|\mathcal{A}|)} + \log(|\mathcal{A}|) \tag{17.3}$$

If we do not know T, we can use the standard doubling trick and the asymptotic bound of (17.3) remains the same. In practice, the $\beta$ parameter can be tuned for the game and can significantly improve performance. Furthermore, there are parameter-free variants e.g. NormalHedge (Chaudhuri et al., 2009).

## 17.2.4 Regret Matching

One particularly simple algorithm is regret matching. One simply plays an action proportionally to the positive regret of that action (Blackwell et al., 1956; Hart and Mas-Colell, 2000). Let $R(a)^+ = \max(R(a), 0)$, regret matching then sets the action probability as follows:

$$\pi_{rm}^t(a) = \frac{R(a)^+}{\sum_{a' \in \mathcal{A}} R(a')^+} \tag{17.4}$$

When $\sum_{a' \in \mathcal{A}} R(a')^+ = 0$, we simply select a uniform random policy $\pi^t_{rm}(a) = \frac{1}{|\mathcal{A}|}$. Note that from the equation 17.4, we see that unlike Hedge (and many others), regret matching is not parametric. This is a powerful property, since in practice we don't have to find parameters that work well for our particular problem. Regret matching enjoys the following bound (17.5).

$$R^T(\pi_{rm}) \leq \sqrt{|\mathcal{A}|T} \tag{17.5}$$

# 17.3   Connection To Nash Equilibrium

We are now ready to show the striking connection from regret to Nash equilibria. Namely, we will connect the average regret $\frac{R^T(\pi)}{T}$ to $\epsilon$-Nash equilibrium. The idea is to use regret minimization in self-play, where we iteratively update policies of the players. Self-play utilities are then used to define the reward vector, which in turn forms the regrets to be minimized (Figure 17.1 and Algorithm 5 ).

At time $t$, a player (regret minimizer) produces their strategy $\pi^t_i$ and the reward vector is computed as the utility against a strategy of the opponent $\pi^t_{-i}$. For a matrix game, the reward vector of the player $i$ is $x^t_i = A\pi^t_{-i}$ and the current regret is: $r^t_i = A\pi^t_{-i} - \pi^t_i A\pi^t_{-i}$ .

## 17.3.1   Convergence Analysis

There is a crucial connection between the average regret $\frac{R^T}{T}$ and the optimality of the average strategy of the self-play. As the regret essentially measures by how much we could have improved and the regrets in self-play are a function of opponent's policy, there is a direct connection to $\epsilon$-Nash equilibrium (Theorem 22).

**Theorem 22.** *Let the accumulated regret in self-play of the players be $R^T_1$ and $R^T_2$. The averaged strategy profile $(\pi_1, \pi_2)$ is then $\epsilon$-Nash for $\epsilon = \frac{R^T_1 + R^T_2}{T}$.*

*Proof.* WLOG we show that player 1 cannot improve by more than $R^T_1 + R^T_2$ by switching to an arbitrary strategy. First, notice that $\sum_t \bar{\pi}_1 A \pi^t_2 = \sum_t \bar{\pi}_1 A \bar{\pi}_2 = \sum_t \pi^t_1 A \bar{\pi}_2$.

$$\max_{\pi^*_1}(\sum_t \pi^*_1(A)\bar{\pi}_2) - \sum_t \pi^t_1(A)\pi^t_2 = \max_{\pi^*_1}(\sum_t \pi^*_1(A)\pi^t_2) - \sum_t \pi^t_1(A)\pi^t_2 \leq R^T_1$$

$$\sum_t \bar{\pi}_1(-A)\bar{\pi}_2 - \sum_t \pi^t_1(-A)\pi^t_2 \leq \max_{\pi^*_2}(\sum_t \bar{\pi}_1(-A)\pi^*_2) - \sum_t \pi^t_1(-A)\pi^t_2 \leq R^T_2$$

$$\sum_t \bar{\pi}_1(A)\bar{\pi}_2 - \sum_t \pi^t_1(A)\pi^t_2 \geq -R^T_2$$

$$\max_{\pi^*_1}(\sum_t \pi^*_1(A)\bar{\pi}_2) - \sum_t \bar{\pi}_1(A)\bar{\pi}_2 \leq R^T_1 + R^T_2$$

$$\max_{\pi^*_1} \pi^*_1(A)\bar{\pi}_2 - \bar{\pi}_1(A)\bar{\pi}_2 \leq \frac{R^T_1 + R^T_2}{T}$$

$\square$

**Algorithm 5** Self-play Regret Minimization in Matrix Games

---

1: **function** RegretMatching(regrets)
2:     positive_regrets = regrets$^+$
3:     pos_regret_sum = $\sum_i$ positive_regrets$_i$
4:     **if** pos_regret_sum = 0 **then return** uniform_policy
5:     **return** positive_regrets / pos_regret_sum

6:
7: **function** RunSelfplay($A \in \mathcal{R}^{n \times n}$, steps)
8:                                                    ▷ Initialize regrets and average policies
9:     $r_1 = \vec{0} \in \mathcal{R}^{1 \times m}$
10:    $r_2 = \vec{0} \in \mathcal{R}^{1 \times n}$
11:    $\overline{\pi}_1 = \vec{0} \in \mathcal{R}^{1 \times m}$
12:    $\overline{\pi}_2 = \vec{0} \in \mathcal{R}^{1 \times n}$
13:    **for** t **do** in steps
14:                                                    ▷ Compute the current policy
15:        $\pi_1 = \text{RegretMatching}(r_1)$
16:        $\pi_2 = \text{RegretMatching}(r_2)$
17:                                                    ▷ Compute the current reward vector
18:        $x_1 = (A\pi_2^\top)^\top$
19:        $x_2 = \pi_1 A$
20:                                                    ▷ Compute the received value
21:        $v_1 = \pi_1 x_1^\top$
22:        $v_2 = \pi_2 x_2^\top$
23:                                    ▷ Compute current regret and add it to the cumulative regret
24:        $r_1 + = x_1 - v_1$
25:        $r_2 + = x_2 - v_2$
26:                                                    ▷ Cumulate average policies
27:        $\overline{\pi}_1 + = \pi_1$
28:        $\overline{\pi}_2 + = \pi_2$
29:                        ▷ Return average policies **return** $\{\overline{\pi}_1/steps, \overline{\pi}_2/steps\}$

---

Theorem 22 connects the average regret and the quality of the average strategy. Thus, all we need is a regret minimizer with a sub-linear regret growth (Hannan consistent) and we have a method that provably converges to Nash equilibria.

### 17.3.2 Possible Bounds

Regret matching and many other regret minimizers provably enjoy $\frac{1}{\sqrt{T}}$ convergence speed of the average regret. This matches the lower bound of the general adversarial case — for any algorithm, it is possible to construct an adversarial sequence where the average regret converges no faster (Nisan et al., 2007).

When the environment is not strictly adversarial, it is possible to leverage the structure in the predictive regret framework. A particularly useful structure is when the reward vector does not drastically change over time — gradual variations (Chiang et al., 2012). The idea of gradual variations can be leveraged in self-play, as we get to "controll" both sides. The reward vector is then unlikely to shift drastically as the agents' policies change only gradually. This structure is then used in the predictive regret framework, where we predict the next reward vector. The better the prediction, the faster the convergence and there are predictive regret algorithms with $\frac{1}{T}$ convergence speed (Syrgkanis et al., 2015; Farina et al., 2019b, 2020). This speed is then optimal for the self-play settings, as it matches the corresponding lower bound (Daskalakis et al., 2011) .

## 17.4 Alternating Updates

When we applied regret minimization in self-play, we updated both players at the same time. The players used regret minimizer to produce their policy $(\pi_1^t, \pi_2^t)$ which is then used to compute the current regret. This is referred to as simultaneous updates (Figure 17.1).

There is another option — alternating updates — where we update players one by one. We use $(\pi_1^t, \pi_2^t)$ to run a regret minimization step for one agent, producing $\pi^{t+1}$. That policy is then used for a regret minimizaton step of the other agent, where the current regret is based on $(\pi_1^{t+1}, \pi_2^t)$ (Figure 17.2).

For simultaneous updates, the regret updates of both players are based on the $(\pi_1^t, \pi_2^t)$ sequence. For alternating updates, one player uses $(\pi_1^t, \pi_2^t)$ while the other one $(\pi_1^{t+1}, \pi_2^t)$. Alternating updates in self-play is thus not symmetric — the reward vectors are not being computed symmetrically as in simultaneous updates. An important caveat is that Theorem 22 holds only for simultaneous updates as black box regret minimization with alternating updates might fail to converge (Farina et al., 2019a). One must use additional assumptions about the regret minimizer to recover the convergence guarantess. Fortunately, regret matching is one such regret minimizer (Burch et al., 2019).

## 17.5 Counterfactual Regret Minimization

In matrix games, we computed the regret against the best action in retrospect. In other words, we compared to all possible pure strategies. But what does that

Figure 17.1: Regret minimization in self-play — simultaneous updates. Both players update their policy at the same time and the regret update is symmetrical as it uses the same strategy pair for both player updates.



Figure 17.2: Regret minimization in self-play — alternating updates. Players update their policy one by one and the regret update of one player uses different strategy pair than the opponent does.

mean in sequential settings? As we know from Section 13.2.3, we can always convert to a matrix game. While possible, this is extremely impractical due to the exponential size explosion.

Fortunately, it is possible to minimize regret directly in the sequential representation of the game. Counterfactual regret minimization (CFR) decomposes the full regret into small individual regrets in each information state (counterfactual regrets). One can then bound the full regret by a sum of all these partial regrets, allowing one to minimize these partial regrets independently (Zinkevich et al., 2007). It is then possible to generalize the idea from sequential games to a more general case of convex sequential decision making (Farina et al., 2018).

### 17.5.1 Counterfactual Values

We first define state and state-action values. Those terms (analogous to their reinforcement learning counterparts) will allow for a particularly easy view and definition of CFR.

The state value $v_i^\pi(h)$ is the expected future reward under policy $\pi$ to player $i$ given the history $h$. The state-action value $q_i^\pi(h, a_i)$ is defined analogously, except that the player $i$ first takes the action $a_i$. Counterfactual reach $P_{-i}^\pi(h)$ of a given history $h$ for player $i$ is the reach probability of that history when player $i$ attempts to reach $h$.

$$P_{-i}^\pi(h) := P_{\mathcal{T}}(h) \prod_{j \in \mathcal{N} \setminus \{i\}} P_j^\pi(h)$$

Now we can define counterfactual-weighted state and state-action values. The counterfactual-weighted state-action value for player $i$ of action $a \in \mathcal{A}_i$ at state $s \in \mathcal{S}_i$ is (17.6).

$$q_{i,c}^{\pi}(s,a) := \sum_{h \in \mathcal{H}(s)} P_{-i}^{\pi}(h) q_i^{\pi}(h,a) \tag{17.6}$$

Having these counterfactually-weighted action-values, we define the corresponding state-values as (17.7).

$$v_{i,c}^{\pi}(s) := \sum_{a \in \mathcal{A}_i(s)} \pi_i(s,a) q_{i,c}^{\pi}(s,a) \tag{17.7}$$

Note that this value is not conditional on reaching $s$, as is standard in reinforcement learning, but instead depends on the counterfactual probability $P_{-i}^{\pi}$ of reaching $s$.

### 17.5.2 Counterfactual Regret

Given a strategy sequence $\pi^0, \ldots, \pi^{t-1}$, we can use counterfactual state and state-action values to define the counterfactual regrets as follows:

$$R_i^t(s,a) = \sum_{k=0}^{t-1} \left( q_{i,c}^{\pi^k}(s,a) - v_{i,c}^{\pi^k}(s) \right) \tag{17.8}$$

$$R_i^t(s) = \max_{a \in \mathcal{A}(s)} R_i^t(s,a) \tag{17.9}$$

Setting $R_i^t(s)^+ = \max(R_i^t(s), 0)$, the key result is that one can bound the full regret by the sum of the positive counterfactual regrets (Theorem 23) (Zinkevich et al., 2007).

**Theorem 23.** *(Zinkevich et al., 2007, Theorem 3)*

$$R_i^t \leq \sum_{s \in \mathcal{S}_i} R_i^t(s)^+$$

Theorem 23 has critical consequences. It decomposes the full regret into a sum of immediate regrets and allows us to independently minimize the immediate state regrets 17.9 (e.g. using regret matching). Furthermore, computing the counterfactual values, regrets and updates can be done via a single tree traversal (Section 17.5.5). The same decomposition idea to individual partial regrets can then be generalized to sequential decision processes for convex sets (Farina et al., 2019a).

### 17.5.3 Average Policy

Combining Theorem 23 and Theorem 22 guarantees that if we use Hannan consistent regret minimizer in each infostate, the average strategy converges to optimal. During the averaging, one must not forget to properly weight the behavioral strategy by the current reach probability (Section 3.1.4).

### 17.5.4 Counterfactual Best Response

Policy $\pi_i$ with non-negative counterfactual regret sum $\sum_{s \in \mathcal{S}_i} R_i(s)^+ = 0$ is optimal and thus also a best response $\pi_i \in \mathbb{BR}$ (Lemma 3). But not every best response $\pi_i'$ has zero regrets. If a state $s \in \mathcal{S}_i$ has zero-reach probability $P_{\pi_i'}(s) = 0$, the policy might form a best response while $R_i(s)^+ > 0$. We thus define a refinement of best response policy — the counterfactual best response. This concept can be thought of as a particular notion of sub-game perfect policy and will be beneficial in later sections.

**Definition 24** (Counterfactual Best Response). *Counterfactual best response against a policy $\pi_i$ is a policy of the opponent with zero counterfactual regrets for all $s \in \mathcal{S}_{-i}$.*

We denote the set of such counterfactual best response policies as $\mathbb{CFBR}(\pi_i)$ and we use $CBRV_{-i}^{\pi_i}(s)$ to denote the counterfactual best response value of a state $s \in \mathcal{S}_2$ — counterfactual value given that the player $i$ follows $\pi_i$ and the opponent follows $\pi_{-i} \in \mathbb{CFBR}(\pi_i)$.

### 17.5.5 Public Tree Implementation



function ComputeValues($s_{pub} \in \mathcal{S}_{pub}$, $d_1 \in \Delta\mathcal{S}_1(s_{pub})$, $d_2 \in \Delta\mathcal{S}_2(s_{pub})$)

Figure 17.3: Public tree implementation of CFR traverses the public tree and send the reach probabilities down the tree and counterfactual values up the tree.

The CFR algorithm iteratively i) produces policy $\pi_i^t(s)$ using counterfactual regret in individual states ii) computes the counterfactual values and updates the regrets iii) averages the policy. One such iteration can be efficiently implemented via a single game tree traversal, where the traversal sends the reach probabilities down the tree (to be used as the $P_{-i}^\pi$ reach weights as in Equation 17.6) and values up the tee (propagating the values as in Equation 17.6). Algorithm 6 then shows an implementation where we traverse a public tree. While there is no strong reason for this choice now and one could traverse the history or infoset trees, a public tree implementation will be beneficial once we will use value functions within the CFR algorithm. Note that the public tree recursive traversal operates with reach probabilities and counterfactual values for all the states of a public state (Figure 17.3), which will become particularly useful for the use of value functions.

**Algorithm 6** Public Tree CFR

---

1: **function** PUBLICTREECFR
2:      **for** $i : [1..T]$ **do**
3:          COMPUTENEXTPOLICIES
4:          UPDATEAVERAGEPOLICIES
5:          COMPUTEVALUES($root_{pub}, 1, 1$)
6:          UPDATEREGRETS
7:      NORMALIZEAVERAGEPOLICIES
8:
9:    ▷ Recursively traverses the public tree, sending reach probabilities down and counterfactual values up.
10: **function** COMPUTEVALUES($s_{pub} \in \mathcal{S}_{pub},\ d_1 \in \Delta\mathcal{S}_1(s_{pub}),\ d_2 \in \Delta\mathcal{S}_2(s_{pub})$)
11:                                     ▷ Terminal public state.
12:      **if** $s_{pub} \in \mathcal{Z}_{pub}$ **then**
13:          **for** $i \in \mathcal{N}$ **do**
14:              **for** $s \in \mathcal{S}_i(s_{pub})$ **do**
15:                  $P_{-i}(h) \leftarrow d_{-i}(s)P_{\mathcal{T}}(h)$
16:                  $v_{i,c}(s) = \sum_{h \in \mathcal{H}(s)} P_{-i}(h)R_i(h)$
         **return**
17:                               ▷ Update the distribution for children.
18:      **for** $i \in \mathcal{N}$ **do**
19:          **for** $s \in \mathcal{S}_i(s_{pub})$ **do**
20:              **for** $sao \in \mathcal{S}_i$ **do**
21:                  $d_i(sao) = d_i(s)\pi_i(a)$
22:                                       ▷ Recursion.
23:      **for** $s_{pub}o_{pub} : \mathcal{S}_{pub}$ **do**
24:          COMPUTEVALUES($s_{pub}o_{pub},\ d_1(s_{pub}o_{pub})$)
25:                            ▷ Compute state-action values.
26:      **for** $i \in \mathcal{N}$ **do**
27:          **for** $s \in \mathcal{S}_i(s_{pub})$ **do**
28:              $v_{i,c}(s) = \sum_{sao \in \mathcal{S}_i} \pi(s,a)v_{i,c}(sao)$
29:              **for** $a \in \mathcal{A}_i(s)$ **do**
30:                  $q_{i,c}(s,a) = \sum_{sao \in \mathcal{S}_i} v_{i,c}(sao)$
31:
32: **function** UPDATEREGRETS
33:      **for** $s \in \mathcal{S}$ **do**
34:          **for** $a \in \mathcal{A}(s)$ **do**
35:              $R(s,a) {+}= q_{i,c}(s,a) - v_{i,c}(s)$
36:
37: **function** COMPUTENEXTPOLICIES
38:      **for** $s \in \mathcal{S}$ **do**
39:          $\pi(s) = $ REGRETMATCHING($R(s)$)

---

## 17.6   CFR+

CFR+ is an algorithm with an interesting history and impressive empirical performance. Initially discovered by a hobbyist scientist (Tammelin, 2014), CFR+ drastically improved state of the art methods and has been used to solve limit Texas hold'em poker — the largest imperfect information game played by humans to be solved to this day (Bowling et al., 2015), and the resulting agent Cepheus has also been open sourced[2].

CFR+ modifies vanilla CFR[3] in three ways: i) it uses regret matching plus ii) it uses alternating updates iii) and it uses linear averaging of the policy.

**i) Regret Matching Plus**   Regret matching plus (RM$^+$) selects actions just like regret matching, except that it caps all the regrets to zero after each iteration. That is, we set $R(a) \leftarrow R(a)^+$. RM$^+$ enjoys the same regret bounds as RM does.

**ii) Alternating updates**   See Section 17.4

**iii) Linear averaging**   The average strategy in Theorem 22 is a uniformly weighted average of the individual (current) policies $\overline{\pi_i} = \frac{1}{T} \sum_{t=1}^{T} \pi_i^t$. CFR+ rather uses linearly increasing weight for the later iterations $\overline{\pi_i} = \frac{2}{(T)(T+1)} \sum_{t=1}^{T} t\pi_i^t$. This averaging has been proven to be sound for CFR+, but does not work in general as the analysis does not apply to CFR.

### 17.6.1   Empirical Performance

CFR+ has some exciting empirical performance. It often greatly outperforms CFR and even many algorithms with a $\frac{1}{T}$ convergence bound. Figure 17.4 compares the convergence speed of CFR+ and CFR on Leduc poker and Glasses.



(a) Leduc poker.          (b) Glasses.

Figure 17.4: Convergence of CFR and CFR+ tabular solvers.

Why does CFR+ drastically outperform CFR and often converges much faster than its theoretical bound of $\frac{1}{\sqrt{T}}$? The correct answer is that no-one really knows, though there are some appealing properties. CFR+ has been initially analysed in (Tammelin et al., 2015), where the authors show that the algorithm also guarantees a sub-linear tracking regret (Herbster and Warmuth, 1998). Furthermore,

---

[2]http://poker.srv.ualberta.ca
[3]As originally introduced in (Zinkevich et al., 2007).

alternating updates in the case of CFR+ are also soundly motivated (Burch et al., 2019).

Unfortunately, it seems that the strong empirical performance of the algorithm is indeed just empirical. While CFR+ often converges much faster than its bound would suggest, there are known cases where the convergence speed is much closer to the $\frac{1}{\sqrt{T}}$ bound, performing no better than CFR. The common belief is that this algorithm's convergence rate is not $\frac{1}{T}$ in the worst case (Burch, 2018; Farina et al., 2019b).

## 17.7 Monte Carlo CFR

While CFR methods traverse the full game tree at each iteration, Monte Carlo methods sample a subset of all such trajectories. Technically, Monte Carlo CFR (MCCFR) is a family of algorithms defined by a partition $\mathcal{Q}$ of the terminal histories $\mathcal{Z}$, sampling a block $\mathcal{Q}_i \in \mathcal{Q}$ at each iterations and updating only the relevant information states for that block. In combination with a proper importance sampling correction term, one can then guarantee that the values and updates are in expectation the same as in CFR, and derive a probabilistic bound for the regrets (Lanctot et al., 2009; Lanctot, 2013).

This has a clear benefit of much faster per-iteration time, but the price one has to pay is that the values are noisy due to the variance introduced by the sampling. Smaller blocks $\mathcal{Q}_i$ provide faster iterations and higher variance. Different sampling schemes thus essentially balance per-speed iteration and variance.

### 17.7.1 Outcome Sampling

We describe a particular member of the MCCFR family that samples a single history at each iteration. That is, the partition $\mathcal{Q} = \mathcal{Z}$, making the notation particularly concise. Furthermore, sampling a single history emphasizes similarities to some reinforcement learning algorithms and will simplify the reduced-variance variant of the algorithm presented in the next section.

Outcome sampling samples a single history $z \in \mathcal{Z}$, updating only the visited information states on the sampled trajectory. To make sure each terminal history is sampled with non-zero probability, we restrict the sampling policy $\xi \in \Delta A(s)$ for all $s$ so that $\xi(s, a) > 0$ everywhere. On each iteration $z \sim \xi$, the algorithm replaces counterfactual value with sampled counterfactual value $\tilde{v}_{i,c}^{\pi}(s|z)$.

$$\tilde{v}_{i,c}^{\pi}(s|z) = \tilde{v}_{i,c}^{\pi}(h|z) = \frac{P_{-i}^{\pi}(h)P^{\pi}(h,z)R_i(z)}{P^{\xi}(z)}, \text{ for } h \in s, h \sqsubseteq z, \qquad (17.10)$$

and 0 for histories that were not played ($h' \not\sqsubseteq z$). Note that the term is simply a counterfactual value computed from a single history $z$ with the importance sampling term $P^{\xi}(z)$. Such estimate is then provably unbiased $\mathbb{E}_{z \sim \xi}[\tilde{v}_{i,c}^{\pi}(s|z)] = v_{i,c}^{\pi}(s, I)$ (Lanctot et al., 2009, Lemma 1). As a result, we can replace $v_{i,c}$ and $q_{i,c}$ in Equation 17.8 with $\tilde{v}_{i,c}$ and $\tilde{q}_{i,c}$ respectively, resulting in the accumulate estimated regrets.

$$\tilde{R}^t(s,a) = \tilde{q}_{i,c}^{\pi}(s,a) - \tilde{v}_{i,c}^{\pi}(s) \qquad (17.11)$$

The corresponding regret bound then requires an additional term $\frac{1}{\min_{z \in \mathcal{Z}} P(z)}$, which is exponential in the length of $z$ (Theorem 24). Similar observations have been made in reinforcement learning (Arjona-Medina et al., 2018).

**Theorem 24.** *(Lanctot et al., 2009, Theorem 5) For any $p \in (0, 1]$, when using outcome-sampling MCCFR where $\forall z \in \mathcal{Z}$ either $P^\pi_{-i}(z) = 0$ or $P^\xi(z) \geq \delta > 0$ at every timestap, with probability $1 - p$, average overall regret is bounded by*
$$R_i^T \leq (1 + \frac{\sqrt{2}}{\sqrt{p}})(\frac{1}{\delta})\Delta_R \frac{\sqrt{|A_i|}}{\sqrt{T}}$$

**Recurrent Computation**

A particularly insightful formulation is to first recurrently formulate sampled state and state-action values $\tilde{q}_i^\pi(h, a|z)$ and $\tilde{v}_i^\pi(h|z)$, just like one would compute and propagate the values in a tree traversal.

$$\tilde{q}_i^\pi(h, a|z) = \begin{cases} \tilde{v}_i^\pi(ha|z)/\xi(h, a) & \text{if } ha \sqsubseteq z \\ 0 & \text{if } h \sqsubset z,\ ha \not\sqsubseteq z \\ 0 & \text{otherwise} \end{cases} \qquad (17.12)$$

$$\tilde{v}_i^\pi(h|z) = \begin{cases} R_i(h) & \text{if } h = z \\ \sum_a \pi(h, a)\tilde{q}_i^\pi(h, a|z) & \text{if } h \sqsubset z \\ 0 & \text{otherwise} \end{cases} \qquad (17.13)$$

These values are then turned into corresponding sampled counterfactual values by

$$\tilde{q}_{i,c}^\pi(h, a|z) = \frac{P^\pi_{-i}(h)}{P^\xi(h)}\tilde{q}_i^\pi(h, a|z) \qquad (17.14)$$

## 17.7.2  Convergence Speed

Unfortunately, MCCFR methods usually converge slower than their full-tree traversal counterparts (even when the comparison takes into account the faster iteration time). The main problem with the sampling variants is that they introduce variance that can have a significant effect on long-term convergence (Gibson et al., 2012). Furthermore, CFR+ loses its magic and MCCFR+ is usually outperformed by MCCFR (Burch, 2018). One way to look at this problem is that it is unlikely to get convergence near $\frac{1}{T}$ in the sampling settings, as the central limit theorem tells us that the average value converges to the expected one at the rate of $\frac{\sigma}{\sqrt{T}}$. Thus one would have to drastically decrease the sampling variance to get rid of this term.

# 17.8  Variance Reduced MCCFR

Variance Reduced MCCFR (VR-MCCFR) (Schmid et al., 2019) and its later variants (Davis et al., 2020) are powerful methods that decrease the per-sample

Figure 17.5: High-level overview of Variance Reduction MCCFR (VR-MCCFR) and related methods. a) CFR traverses the entire tree on every iteration. b) MCCFR samples trajectories and computes the values only for the sampled actions, while the off-trajectory actions are treated as zero valued. While MCCFR uses importance sampling weight to ensure the values are unbiased, the sampling introduces high variance. c) VR-MCCFR follows the same sampling framework as MCCFR, but uses baseline values for both sampled actions (in blue) as well as the off-trajectory actions (in red). These baselines use control variates and send up bootstrapped estimates to decrease the per-iteration variance thus speeding up the convergence.

variance in MCCFR, resulting in orders of magnitude faster convergence. VR-MCCFR is also another contribution of this thesis.

At the core of the method is the general idea of control variates. The idea is to estimate the value of non-sampled actions and use the control variates technique to make sure the expectation of the action-values remains unbiased regardless of the estimates. While any estimate leads to unbiased values, poor estimates can increase the variance rather than decrease it. The closer the estimates are to the true value, the lower the variance. Furthermore, VR-MCCFR recursively combines the estimates as it propagates values up the tree, propagating the benefits of the value estimates. See Figure 17.5 for comparison of CFR, MCCFR and VR-MCCFR.

## 17.8.1 Control Variates

Suppose we are trying to estimate a mean from samples $X = (X_1, X_2, \ldots, X_n)$. The Monte Carlo estimator is then simply $\tilde{X}^{mc} = \frac{1}{n} \sum_{i=1}^{n} X_i$. A control variate is a random variable $Y$ with a known mean $\mathbb{E}[Y]$ to be paired with the original variable, resulting in a new random variable $Z_i = X_i + c(Y_i - \mathbb{E}[Y])$ (with a corresponding Monte Carlo estimator $\tilde{Z}^{mc}$). Since $E[Z_i] = \mathbb{E}[X_i]$ for any $c$ we can use $\tilde{Z}^{mc}$ instead of $\tilde{X}^{mc}$ with variance $Var[Z_i] = Var[X_i] + c^2 Var[Y_i] + 2c Cov[X_i, Y_i]$ (Owen, 2016). So when $X$ and $Y$ are positively correlated and $c < 0$, variance is reduced when $Cov[X, Y] > \frac{c^2}{2} Var[Y]$.

## 17.8.2 Baseline Enhanced Estimated Values

VR-MCCFR constructs value estimates using control variates. We first define an action-dependent baseline $b_i(s, a)$ to be used as a control variate to approximate or be correlated with $\mathbb{E}[\tilde{q}_i^\pi(s, a)]$. We also define a sampled baseline $\tilde{b}_i(s, a)$ for which $\mathbb{E}[\tilde{b}_i(s, a)] = b_i(s, a)$. We can now readily use these terms in the control variate framework to construct a new baseline-enhanced estimate for the state-action values:

$$\tilde{q}_{i,c}^{b,\pi}(s,a) = \tilde{q}_{i,c}^{\pi}(s,a) - \tilde{b}_i(s,a) + b_i(s,a) \tag{17.15}$$

First, note that $\tilde{b}_i$ is a control variate with $c = -1$ and thus the expectation remains unchanged (Lemma 25). Equation 25 essentially corresponds to a local adjustment of sampled counterfactual values in a single state. Figure 17.6 illustrates such local computation for CFR, outcome sampling MCCFR and outcome sampling VR-MCCFR.

**Lemma 25.** *(Schmid et al., 2019, Lemma 1) For any $i \in \mathcal{N} - \{c\}, \pi_i, s \in \mathcal{S}_i, a \in \mathcal{A}_i(s)$, if $\mathbb{E}[\tilde{b}_i(s,a)] = b_i(s,a)$ and $\mathbb{E}[\tilde{q}_{i,c}^{\pi}(s,a)] = q_{i,c}^{\pi}(s,a)$, then $\mathbb{E}[\tilde{q}_{i,c}^{b,\pi}(s,a)] = q_{i,c}^{\pi}(s,a)$.*



Figure 17.6: a) CFR computes the exact value by traversing all the actions. b) Outcome sampling MCCCFR samples a single action and uses zero for the values of the non-sampled actions. c) VR-MCCFR uses baselines to adjust the estimates of all the actions.

**Recursive Bootstrapping**

Given the recursive nature of the counterfactual values propagation, we can propagate the already baseline-enhanced counterfactual values rather than the noisy sampled values. This allows us to propagate the benefits up the tree rather than using the control variates trick in isolation. Just as we recursively formulated sampled state and state-action values (Equation 17.12 and 17.12) that were then used to formulate the sampled counterfactual values for outcome sampling, we recursively formulate the baseline-enhanced variants.

$$\tilde{q}_i^{b,\pi}(h,a|z) = \begin{cases} b_i(s(h),a) + \frac{\tilde{v}_i^{b,\pi}(ha|z) - b_i(s(h),a)}{\xi(h,a)} & \text{if } ha \sqsubseteq z \\ b_i(s(h),a) & \text{if } h \sqsubset z, ha \not\sqsubseteq z \\ 0 & \text{otherwise} \end{cases} \tag{17.16}$$

and

$$\tilde{v}_i^{b,\pi}(h|z) = \begin{cases} u_i(h) & \text{if } h = z \\ \sum_a \pi(h,a)\tilde{q}_i^{b,\pi}(h,a|z) & \text{if } h \sqsubset z \\ 0 & \text{otherwise} \end{cases} \tag{17.17}$$

Just like in outcome sampling, we then turn these estimates of state and state-action values into estimates of counterfactual values, which are in turn used to update the regrets.

$$\tilde{q}_{i,c}^{b,\pi}(s,a|z) = \tilde{q}_{i,c}^{b,\pi}(h,a|z) = \frac{P_{-i}^{\pi}(h)}{P^{\xi}(h)} \tilde{q}_i^{b,\pi}(h,a|z) \qquad (17.18)$$

The resulting values can be then Note that these values collapse to the original outcome sampling for $b_i(s,a) = 0$ ($\tilde{q}_i^{b,\pi}(s,a)$ becomes $\tilde{q}_i^{\pi}(s,a)$). Outcome sampling (and other MCCFR variants) can thus be viewed as VR-MCCFR algorithm with particular choice of zero baseline.

**Lemma 26.** *(Schmid et al., 2019, Lemma 2) Let $\tilde{q}_{i,c}^{b,\pi}$ be defined as in Equation 17.18. Then, for any $i \in \mathcal{N} - \{c\}, \pi_i, s \in \mathcal{S}_i, a \in \mathcal{A}(s)$, it holds that $\mathbb{E}_z[\tilde{q}_{i,c}^{b,\pi}(s,a|z)] = q_{i,c}^{\pi}(s,a)$.*

### 17.8.3 Choice of Baselines

While any choice of baseline leads to unbiased estimates, the closer the estimate to the true value, the better (Theorem 27). A typical choice of the estimates is to simply use values from previous iterations (e.g. with exponentially decaying weight), as this value is unlikely to drastically change. This simple choice already leads to a drastic speed improvement over MCCFR.

**Theorem 27.** *(Gibson et al., 2012, Theorem 2) For some unbiased estimator of the counterfactual values $\tilde{v}_i$ and a bound on the difference in its value $\tilde{\Delta}_i = |\tilde{v}_i(\pi, I, a) - \tilde{v}_i^{\pi}(I, a')|$, with probability $1 - p$*

$$\frac{R_i^T}{T} \le \left( \tilde{\Delta}_i + \frac{\sqrt{\max_{t,s,a} Var[r_i^t(s,a) - \tilde{r}_i^t(s,a)]}}{\sqrt{p}} \right) \frac{|\mathcal{S}_i||\mathcal{A}_i|}{\sqrt{T}}.$$

Another key result is that there exists a perfect baseline that leads to zero-variance estimates at the updated information sets (Lemma 28). Since the oracle baseline corresponds to the exact value, it is as hard to compute as the full tree traversal of CFR and thus impractical. On the other hand, it can be used to analyze the performance of the algorithm and how our choice of baseline compares to the optimal one. Note that even with the oracle baseline, VR-MCCFR is still not identical to CFR as VR-MCCFR only updates a subset of information states (albeit with the exact counterfactual values) while CFR updates the entire tree.

**Lemma 28.** *(Schmid et al., 2019, Lemma 3) There exists a perfect baseline $b^*$ and optimal unbiased estimator $\tilde{v}_i^{*,\pi}(s,a)$ such that under a specific update scheme: $Var_{h,z \sim \xi, s \in \mathcal{S}, h \sqsubseteq z}[\tilde{v}_i^{*,\pi}(h,a|z)] = 0$.*

### 17.8.4 Convergence Speed and Variance

We evaluate the algorithm on Leduc poker with exponentially decaying weight for the average baseline (weight $\alpha = 0.5$), comparing MCCFR, MCCFR+, VR-MCCFR, VR-MCCFR+, and VR-MCCFR+ with the oracle baseline. Variance reduced variants converge significantly faster than MCCFR and the relative speedup grows as the baseline improves during the computation. VR-MCCFR+ achieves two orders of magninute speedup compared to MCCFR. VR-MCCFR+

Figure 17.7: Convergence of exploitability for different MCCFR variants on logarithmic scale. VR-MCCFR converges substantially faster than plain MCCFR. VR-MCCFR+ bring roughly two orders of magnitude speedup. VR-MCC with oracle baseline (actual true values are used as baselines) is used as a bound for VR-MCCFR's performace to show possible room for improvement. When run for $10^6$ iterations VR-MCCFR+ approaches performance of the oracle version. The ribbons show 5th and 95th percentile over 100 runs.

with the oracle baseline initially outperforms VR-MCCFR+, but they get closer as time progresses and the learned baseline improves (Figure 17.7).

To verify that the observed speedup is indeed thanks to the lower variance, we measured the variance of counterfactual value estimates for MCCFR+ and MCCFR (Figure 17.8). We sampled $1,000$ trajectories rooted in visited information sets, with each trajectory sampling a different estimate of the counterfactual value. While the variance of value estimates in MCCFR seems to be more or less constant during the computation, the variance of VR-MCCFR and VR-MCCFR+ value estimates is substantially lower and continues to decrease.



Figure 17.8: The empirical variance of the counterfactual values is significantly lower for VR-MCCFR and further decreases as the estimates improves over time.

### 17.8.5   Connection to RL Baseline Methods

VR-MCCFR has close connections to well-used baselines in reinforcement learning methods (e.g. REINFORCE algorithm with baselines (Williams, 1992)). That is simply because these baseline techniques of reinforcement learning are simply control variates in disguise. An important difference in imperfect information settings of VR-MCCFR, state-action baselines perform significantly faster than state baselines. In single agent environments of reinforcement learning, state-action baselines tend to perform no better than state baselines (Tucker et al., 2018). This is likely due to the fact that the optimal policy is deterministic and one ends up sampling the action with the best value. In imperfect information, this is not the case as we have to keep sampling all the actions at each iteration to guarantee convergence.

# 18. Offline Solve III - Approximate and Abstraction Methods

Tabular methods are limited to games small enough for strategies to be stored explicitly. In the same way there is no hope to use tabular methods to solve chess, we can not hope to use this approach for large imperfect information games. For perfect information games, a common approach is then to use the online search methods. But as this was for a long time thought to be impossible in imperfect information, there are some other non-search methods one can opt for.

## 18.1   Abstraction Methods

First method is to simply use the tabular methods, but apply those on a smaller, abstracted game (Figure 18.1). The hope is that if the abstracted game strategically resembles the full game, the resulting strategy will approximate the optimal solution.



Figure 18.1: The full game (a) is first abstracted into a smaller, tractable abstracted game (b). The abstracted game is then solved via a tabular offline method, producing a policy (d). In order to play the original game, we need to translate a state from the full game back into the abstracted game. A translation issue arises as there are necessarily fewer states in the abstracted game.

This abstraction approach indeed was for a very long time state of the art method for large imperfect information games (Sandholm, 2010; Johanson, 2016). Considerable research was then devoted to larger and smarter abstraction techniques (Hawkin et al., 2011; Johanson et al., 2013; Ganzfried and Sandholm, 2014; Schmid et al., 2015).

The ideal case is a lossless abstraction (Gilpin and Sandholm, 2007), but this is not possible for large games where one usually solves as large of an abstraction as possible. But larger abstractions do not necessarily lead to a better exploitability in the original game. Even more surprisingly, this might not be the case even if

one abstraction is strictly more finer-grained (Waugh et al., 2009). In practice though, larger abstractions tend to produce stronger policies (Johanson, 2016)

### 18.1.1 Translation Issue

The fundamental problem with abstraction methods is the translation. This became apparent relatively early on in work on no-limit poker, where the agents are allowed to bet any amount of chips up to their stack size. For a stack size of $20,000$ chips, there are about $20,000$ actions corresponding to different bet sizes, making this action tree intractably large. Abstraction can thus only contain a subset of the actions, misunderstanding many of them — e.g. the closest states for the opponent's bet of $12,000$ could be $10,000$ and $15,000$. An agent then has no way to properly understand pot-odds in such states — a key concept for optimal play in poker (Chen and Ankenman, 2006).

Soft translation tries to interpolate the strategy between close states in the abstraction (Schnizlein et al., 2009). But soft translation is in many ways just a Band-Aid on a bullet wound. Local best response has been able to exploit even the best of the abstraction agents by huge margins (Lisy and Bowling, 2017). It is fair to say that abstraction based methods are now mostly obsolete. To see the naivety of the abstraction approach, consider the abstraction approach for a chess agent (Figure 18.2).



Figure 18.2: Abstraction method for chess requires constructing an abstracted game (b), and then to map the states and actions between the original game (a) and the abstracted game (b)

### 18.1.2 Local Best Response

Local best response (LBR) has shown that due to the translation issue, even the best abstraction-based no-limit poker agents are heavily exploitable by relatively simple techniques (Lisy and Bowling, 2017). As exact best response value in the large game of no-limit poker is intractable, local best response approximates the value by approximating the best-responding policy. By fixing the opponent, the environment becomes a single-agent environment where the optimal policy is best-responding (Bowling, 2003) and we can thus use any single-agent algorithm to find (or to approximate) such policy. The performance against the resulting policy than serves as a lower bound on the agent's exploitability.

Local best response uses a small search tree with a poker-heuristic value function. In its search tree, it only looks two actions forward (at most one action of the opponent) and considers only a subset of the aviable actions (e.g. fold, call pot bet). For its value function, it uses an expected rolout value when both players follow the call action until the end of the game (this value can be efficiently and exactly computed for Texas hold'em poker).

Table 18.1 reports the performance of different poker agents against the local best response technique. Interestingly, all the agents are substantially more exploitable than simply folding each hand (it is easy to verify that the exploitability of "always fold" policy is 750$mbb$).

Furthermore, Table 18.1 includes an agent that uses no card abstraction and uses the F,C,P,A action abstraction. When LBR considers only these actions, there is no translation issue and LBR indeed fails to exploit this agent. But once LBR considers more actions, it also exploits this agent and further confirms the translation issue of abstraction techniques.

| LBR Pre-flop Actions | F, C | C | C | C |
|---|---|---|---|---|
| LBR Flop Actions | F, C | C | C | 56bets |
| LBR Turn Actions | F, C | F, C, P, A | 56bets | F, C |
| LBR River Actions | F, C | F, C, P, A | 56bets | F, C |
| Hyperborean (2014) | 721 ± 56 | 3852 ± 141 | 4675 ± 152 | 983 ± 95 |
| Slumbot (2016) | 522 ± 50 | 4020 ± 115 | 3763 ± 104 | 1227 ± 79 |
| Act1 (2016) | 407 ± 47 | 2597 ± 140 | 3302 ± 122 | 847 ± 78 |
| Full Cards [100 BB] | -424 ± 37 | -536 ± 87 | 2403 ± 87 | 1008 ± 68 |

Table 18.1: Local best response performance against strong abstraction-based poker agents in no-limit Texas hold'em poker [mbb/g]. The list of actions considered by the technique varies in different poker rounds, and we report four different configurations. Actions are (F)old, (C)all (P)ot bet and (A)ll-in. For the full list of 56 bets see (Lisy and Bowling, 2017). Note that a policy that simply always folds each hand is exploitable by 750$mbb$.

## 18.2 Non-tabular Methods

Another option for solving games that are intractable for the tabular methods is to represent the strategies and other necessary values (e.g. regrets) implicitly (e.g. using neural networks), relying on the generalization power of the networks. Together with methods that do not have to traverse the full game tree (e.g. MCCFR or VR-MCCFR), this approach allows for algorithms that approximately solve very large games.

The first non-tabular method, RCFR (Waugh et al., 2015) used regression trees and needed to keep track of a dataset that was substantially larger than the full game. This limited the approach from being useful for exactly what one would hope to use non-tabular methods for — large games. The first two successful

instances of non-tabular method in large games are then DeepCFR (Brown et al., 2018a) and double[1] neural counterfactual regret minimization (DNCFR) (Li et al., 2018), both combine deep learning and MCCFR and being published around the same time. Both algorithms use separate networks to represent i) regrets and ii) the average policy. MCCFR is used to collect batches for training, using multiple trajectories to decrease the variance of the training data.

Single DeepCFR follows in these steps and shows how to replace the average policy network by storing a buffer of past regret networks (Steinberger, 2019). The idea is that since the current strategy is a simple function of regrets, storing some of the previous networks amounts to storing these current strategies. One can then approximate the average policy from previous strategies in the buffer. Later work then combines modern VR-MCCFR methods into the DREAM algorithm (Steinberger et al., 2020). ARMAC is another MCCFR method that is closely related to DREAM (Gruslys et al., 2020). Finally, there are recent methods that do not build on the CFR family such as properly modified follow the regularized leader algorithm (Perolat et al., 2020).

---

[1]Not to be confused we double q-learning methods in reinforcement learning where the two networks are used to eliminated bias. Double here refers to separate networks for regrets and policies.

# 19. Online Settings

This chapter re-visits the online setting as introduced in Section 7. There are only minor differences to account for the formalism of the imperfect information games, and most of definitions remain intact. The only notable difference is a particular global consistency connection, where Theorem 10 no longer holds, making search in imperfect information more challenging.

Furthermore, optimal policies for imperfect information games might have to be stochastic (Corollary 17.1). Repeated games then allow us to correctly capture the dynamics of online algorithm that an "averaged" offline policy never could. Consider a simple online algorithm for rock-paper-scissors that simply produces the sequence (rock, paper, scissors, rock, paper, ...). Despite its "average" strategy being optimal, this online algorithm is clearly highly exploitable.

## 19.1 Repeated Game

A repeated game with imperfect information is just like a repeated game with perfect information, except that the individual matches consist of imperfect information games. The definition has to reflect this, and we will stick with the FOSG formalism.

Formally, The repeated game $p$ consists of a finite sequence of $k$ individual matches $p = (z_1, z_2, \ldots, z_k)$, where each match $z_i \in \mathcal{Z}$ is a sequence of world states and actions $z_i = (w_i^0, a_i^0 \, w_i^1, a_i^1 \ldots, a_i^{l_i-1}, w_i^{l_i})$ (ending in a terminal world state $w_i^{l_i}$). For each visited world state in the match, there is a corresponding player state (information state) $s_i(w_i^t)$, i.e. their private perspective of the game. For perfect information games, the notion of player state and world state collapsed as the player gets to observe the world perfectly.

## 19.2 Online Algorithm

Just as in Section 7.3, an online algorithm $\Omega$ then simply maps a player's state observed during a match to a strategy (Definition 11).

**Definition 11.** *(Šustr et al., 2020, Definition 2) Online algorithm $\Omega$ is a function $s \times \Theta \mapsto \Delta(\mathcal{A}_i(s)) \times \Theta$ that maps an information state $s \in \mathcal{S}$ to a strategy $\Delta(\mathcal{A}_i(s))$, while possibly making use of algorithm's state $\theta \in \Theta$ and updating it. We denote the algorithm's initial state as $\theta_0$.*

Given two players $\Omega_1, \Omega_2$, we use $P_{\Omega_1,\Omega_2}^k$ to denote the distribution of gameplays when these two players face each other. The average reward of $i$ is $R(p) = 1/k \sum_{i=1}^{k} u_p(z_i)$ and we denote $\mathbb{E}_{p \sim P_{\Omega_1,\Omega_2}^k}[\mathcal{R}_p]$ to be the expected average reward when the players play $k$ matches.

## 19.3 Soundness

Definition of sound algorithm also remains unchanged

**Definition 12.** *(Šustr et al., 2020, Definition 4) For an $\epsilon$-sound online algorithm $\Omega$, the expected average reward against any opponent is at least as good as if it followed an $\epsilon$-Nash equilibrium fixed strategy $\pi$ for any number of matches $k$:*

$$\forall k \forall \Omega_2 \; : \; \mathbb{E}_{p \sim P^k_{\Omega, \Omega_2}}[R(p)] \geq \mathbb{E}_{p \sim P^k_{\pi, \Omega_2}}[R(p)]. \tag{7.1}$$

## 19.4 Search Consistency

Concepts local, global and strongly global consistencies remain unchanged. An important discrepancy is that Theorem 10 does not hold in imperfect information.

**Theorem 10.** *(Šustr et al., 2020, Theorem 8) In perfect information games, an algorithm that is locally consistent with a subgame perfect equilibrium is sound.*

Local consistency is not sufficient in imperfect even if the algorithm is consistent with sub-game perfect equilibrium. To illustrate this, we will construct a simple game where following is true.

- The game has multiple sub-game perfect Nash equilibria and $\pi^a, \pi^b \in \mathbb{NEQ}$.

- There are two states for the second player $\mathcal{S}_2 = \{s_1, s_2\}$

- Following $\pi^a$, $\pi^b$ in $s_1$ and $s_2$ respectively produces highly exploitable policy.



Figure 19.1: Coordinated matching pennies. First player makes a private action (H)eads or (Tails). Follows a chance player, selecting uniformly one of the publicly observable actions (L)eft or (R)ight. There are now two infostates for the second player $s_1 = [L||], s_1 = [R||]$ and four histories. Second player now chooses (L)eft or (R)ight. In their first infoset $s_1$, they win if they match the opponent's action. In their second infoset $s_2$, they win if they select the opposite action.

Coordinated matching pennies (Figure 19.1) satisfies all three properties. The two information states of player two are $s_1 = [L||], s_1 = [R||]$. It is easy to verify that any optimal policy needs to satisfy $\pi_2(s_1)(H) = \pi_2(s_2)(T)$. In other words, the player has to coordinate between these two information states. Thus

we can construct optimal policies $\pi_1^a, \pi_2^a$ as $\pi_2^a(s_1, L) = 1, \pi_2^a(s_1, L) = 0$ and $\pi_2^b(s_1, L) = 0, \pi_2^a(s_1, L) = 1$. Now if the search technique follows $\pi_1^a$ in $s_1$ and $\pi_2^a$ in $s_2$, the resulting policy is highly exploitable.

While there are multiple notions of sub-game perfect policies in imperfect information settings, note that two information states $s_1, s_2$ have no predecessors and thus this counter-example holds regardless of the notion choice.

### 19.4.1   Online Outcome Sampling



(a) Coordinated matching pennies.

(b) Kuhn poker.

Figure 19.2: Exploitability of individual MCCFR policies and tabularized OOS.

A particularly interesting example of an algorithm that is only locally consistent is the Online Outcome Sampling (OOS) (Lisý et al., 2015). At high level, OOS runs the offline MCCFR algorithm (Section 17.7) in the full game (while also gradually building the tree), parameterized to increase the sampling probability of the current information state. The algorithm then plays based on the resulting strategy for that particular state. The problem is that these individual MCCFR runs can converge to different equilibria as the MCCFR is parameterized differently in each information state. In other words, the OOS algorithm exactly suffers from the fact that it is only locally consistent. Figure 19.2 compares exploitability of the individual MCCFR policies $\pi_2^a(s_1), \pi_2^b(s_2)$ that OOS is locally consistent with, and the exploitability of the tabularized policy the OOS actually follows $\Omega(s_1) = \pi_2^a, \Omega(s_2) = \pi_2^b$. The experiment is run on coordinated matching pennies and Kuhn poker, where the individual MCCFR runs were then parametrized to favor a particular Nash equilibrium — see the Appendix of Šustr et al. (2020) for all the experimental details.

# 20. Search

This chapter introduces online search algorithms for imperfect information. Just as in perfect information, sub-games and value function serve as the key building blocks. And while these concepts are already more complex in imperfect information, the matter is farther complicated by the fact that strong global consistency is not trivial to achieve as we have to coordinate the policy across information states. Furthermore, the base algorithm used for the search techniques is counterfactual regret minimization. This algorithm is particularly suitable for value functions, as it already decomposes the full regret to the sum of individual partial infostate regrets.



Figure 20.1: a) Sub-game (re)solve full lookahead. b) Continual re-solving with full lookahead. c) Continual re-solving with limited lookahead and value functions.

We first start with a single step search algorithm that reasons about the current sub-game (Figure 20.1a), analogically to full lookahead online minimax algorithm of Section 8.1. The full lookahead setting is where we will deal with the complication of consistency. We begin with unsafe re-solving, an unsound approach that motivates the key ingredient in the future sound methods — constrain counterfactual values that ensure the consistency of individual local policies. These values are then used within another key concept — re-solving games (and re-solving gadgets), where a particular gadget game guarantees to produce a sub-game satisfying the constraints. Follows another contribution of the thesis, where a particular modification to the gadget allows to improve (refine) a sub-game policy while still ensuring the strong global consistency.

We then extend this algorithm to multiple steps, where we continually use safe re-solve method during the gameplay — the continual re-solving algorithm (Figure 20.1b). Continual resolving with full lookaheads is analogous to the full looakehad minmax search algorithm. The main complication is that the algorithm keeps updating the constrains to be used for the safe resolving process.

We then combine this algorithm with generalized value functions, allowing to use limited lookaheads within the CFR algorithm. This results in the final algorithm — continual resolving with limited lookaheads and value functions (Figure 20.1c). The final algorithm is the culmination of all the building blocks introduced in this book, presenting a safe search algorithm for imperfect information games.

We finish with the final contribution of the thesis — Deepstack. DeepStack was the first to introduce the continual re-solving algorithm, combining sound search and learned value functions for poker (Moravčík et al., 2017). The combination of continual re-solving and neural networks as value functions resulted in substantial improvement over the previous methods and DeepStack became the first program to beat professional human players in no-limit Texas hold'em poker.

## 20.1 Full Lookahead Unsafe Re-solving

In imperfect information, the notion of the sub-game was defined as a combination of i) public state and ii) distribution over the players' infostates (Section 15.3). We first analyze what happens if we re-solve a sub-game where the distribution over the player's initial states is computed using an optimal policy of the full game. That is we:

1. Compute optimal policy profile $(\pi_i^*, \pi_{-i}^*)$.

2. Select a public state $s_{pub}$ as a root of a sub-game.

3. Compute reach distribution $\Delta(\mathcal{S}_i(s_{pub})), \Delta(\mathcal{S}_{-i}(s_{pub}))$ using the agent's respective policies and the corresponding reach probabilities $P^\pi$.

4. Solve the sub-game defined by $s_{pub}, \Delta(\mathcal{S}_i(s_{pub})), \Delta(\mathcal{S}_{-i}(s_{pub}))$, producing sub-game policy profile $(\pi_i^{sub*}, \pi_{-i}^{sub*})$.

Is the resulting $(\pi_i^{sub*}, \pi_{-i}^{sub*})$ consistent with $(\pi_i^*, \pi_{-i}^*)$? Is such algorithm sound? In perfect information, the answer is positive as the search minmax algorithm fits this framework. And while this approach has has also been used in imperfect information (Ganzfried and Sandholm, 2015) (referred to as the "sub-game solving"), it both loses its theoretical guarantees and empirically produces highly exploitable policies.

Surprisingly, this idea breaks already in rock-paper-scissors (Ganzfried and Sandholm, 2015; Burch et al., 2014). Consider a sub-game rooted in the second public state (just after the first player has acted). The public state contains a single decision infostate of the second player $\{[||]\}$, and three possible states for the first player $\{[|R|], [|P|], [|S|]\}$. As the optimal strategy for the first player is to play uniformly in the first state of the game (prior to reaching this sub-game), distribution over their sub-game states is also uniform. The resulting initial distributions of the sub-game are then $\Delta(\mathcal{S}_1(s_{pub})) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and $\Delta(\mathcal{S}_2(s_{pub})) = (1)$. See Figure 20.2a for illustration.

It is easy to verify that any policy in this sub-game is optimal. But while any policy is optimal in the sub-game, not any strategy is optimal in the original game (e.g. playing always rock is highly exploitable). Thus simply solving a sub-game where the distribution is constructed from prior optimal policies is not even a locally consistent approach. This approach is also referred to as the unsafe re-solving.

Figure 20.2: a) Public sub-game with the initial distribution over the states based on a optimal policy in the full game. b) Counterfactual values corresponding to that optimal policy.

### 20.1.1 Why Did It Break?

After closer inspection of the counterexample, we can investigate why exactly this approach failed. Understanding the issue will allow us to design a new, safe approach that will be crucial in constructing the final algorithms.

Why was the strategy optimal for the sub-game but not locally consistent? Remember that the optimal policy is maximizing the worst case value. When we constructed the sub-game and solved for the optimal policy, we allowed the opponent to best-respond in the sub-game, but we fixed the reach probabilities (sub-game initial distribution) using an optimal policy prior to reaching this sub-game. The opponent could not change their policy prior to this sub-game and alter their reach probabilities.

Why was it beneficial for the opponent to change their policy prior to this sub-game? Consider their sub-game values $V_1$ at the beginning of the sub-game. Under the optimal (uniform) policy of the second player, the values are $V_1 = (0,0,0)$. But after the unsafe re-solving of the sub-game, the values for these states become $(0,1,-1)$. When weighted by the uniform reach probability, the sub-game value remains unchanged $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^\top (0,0,0) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})^\top (0,1,-1) = 0$. The issue is that it is now beneficial for the opponent to change their policy prior to this sub-game, by reaching their second state more often and improve their total best-response value.

**Value Perspective**

Useful viewpoint at the issue is to realize that in perfect information, there is a single root state of the opponent and a single counterfactual value of the opponent that we care to minimize. In imperfect information, there are multiple potential states of the opponent $\mathcal{S}_{-i}(s_{pub})$ and corresponding value vector $V_{-i}(s_{pub})$. What policy should be then preferred, policy $\pi_i^a$ where $V_{-i} = (1,-1)$ or $\pi_i^b$ where $V_{-i} = (-1,2)$? Fixed initial beliefs provide a weighting over those individual values resulting in a single weighted value to be minimized. But as the weights (reach probabilities) can change, so can the resulting value.

Note that if at any time an imperfect information $s_{pub}$ consists of only a single opponent's state, there is no issue and the same can be argued for a sub-game

where there are no prior decisions of the opponent,

## 20.2   Safe Re-solving

Safe re-solving methods make sure the opponent can't increase their sub-game value by altering their reach probabilities. To prevent this, we need to find a policy in the sub-game for which the individual counterfactual values[1] of the opponent remain the same as under the original optimal policy of the full game (or even better, are no greater than the original ones). This then guarantees strong global consistency with a policy having the re-solved values.

Unsafe re-solve used i) public state $s_{pub}$ ii) player's reach $\Delta(\mathcal{S}_i(s_{pub}))$ iii) opponent's reach $\Delta(\mathcal{S}_{-i}(s_{pub}))$ and solved the related sub-game. Safe re-solving on the other hand uses

(i) public state $s_{pub}$.

(ii) player's reach $\Delta(\mathcal{S}_i(s_{pub}))$.

(iii) opponent's counterfactual values $V_{-i}^{bound}(s_{pub})$.

Safe re-solving then produces a policy for which $V_{-i}(s_{pub}) \leq V_{-i}^{bound}(s_{pub})$. Note that this is always possible as the original optimal policy we are re-solving satisfies this constraint.

### 20.2.1   Optimization Formulation

The $V_{-i}(s_{pub}) \leq V_{-i}^{bound}(s_{pub})$ constraints is relatively straightforward to add to the linear programming formulation (Section 16.3.2), as the vector $u$ already present in the optimization corresponds to the negative (counterfactual) best response values of the opponent.

$$
\begin{aligned}
&\min_{u,v} e^\top u \\
&Fy = f, E^\top u - Ay \geq \mathbf{0}, y \geq \mathbf{0} \\
&u_s \leq V_{-i}(s_{pub}, s) \,\forall s \in \mathcal{S}_i(s_{pub})
\end{aligned} \tag{20.1}
$$

Linear optimization problem 20.1 then produces sub-game policy for the player $i$ that is consistent with the original, full game optimal policy we used to derive $\Delta(\mathcal{S}_i(s_{pub}))$ and $V_{-i}^{bound}(s_{pub})$ (Burch et al., 2014; Moravcik et al., 2016).

### 20.2.2   Gadget Formulation

But there is a more elegant and general approach, the gadget game formulation (Burch et al., 2014). The idea is to carefully construct a game for which the resulting optimal policy is guaranteed to satisfy the constraints. Technically, we will construct a gadget game for which there exists a trivial mapping between an optimal solution of the gadget game and the desired policy.

---

[1]Counterfactual best response values.

The core of the construction is that we insert a gadget on top of the sub-game (often referred to as the "CFR-D gadget"). It starts with a chance node that first distributes the opponent into their root information states $\mathcal{S}_{-i}(s_{pub})$, and the opponent then gets to choose (in each respective infostate) whether to (t)erminate and receive the constraint values, or to (f)ollow and play the original sub-game. This clever construction forces the player to play the sub-game so that the opponent's values are no greater than the corresponding (t)erminate values. Figure 20.3 shows a high-level idea of the gadget construction and Figure 20.4 shows a concrete construction for rock-paper-scissors. Note that for the sake of clarity, the figures omit the initial chance state.

The beauty of the gadget game formulation is that it simply constructs a new game of similar size as the sub-game. This allows us to use any algorithm for solving imperfect information games, e.g. any method from the CFR family.



Figure 20.3: a) The gadget game is constructed for a specific public sub-game, and the construction needs i) player's reach probabilities and ii) the opponent's counterfactual values. b) The opponent gets to play first, and in each infoset chooses whether to (t)erminate and receive the constraint values, or to (f)ollow and play the original sub-game.



Figure 20.4: a) Sub-game of rock-paper-scissors. b) CFRD gadget for re-solving the sub-game. Note that the zero values for the (T)erminate action come from the optimal solution.

## Stackelberg Equilibria

The same idea of gadget construction has recently been generalized for the solution concept of Stackelberg equilibria (Ling and Brown, 2021). The paper follows

two steps. First, it shows that safe re-solving for Stackleberg equilibria can be achieved using a carefully constructed upper and lower bound constraints (compared to just upper bound). Second, it introduces modified gadget construction for the lower bound constraints (and uses the "CFRD gadget" construction for the upper bound constraints).

## 20.3   CFR-D

The notion of safe re-solving and the gadget formulation was first introduced by the CFR-D algorithm (Burch et al., 2014). The CFR-D algorithm decomposes the game tree into trunk and sub-games, storing and averaging the policy only for the trunk. During each CFR iteration, individual sub-games (using the reach distribution defined by the current CFR policy) are solved for and strategies discarded, returning only the individual values. The trunk computation thus uses a form of limited lookahead and value function (as described in Section 20.6). During playtime, CFR-D simply uses the stored trunk policy until it reaches one of the sub-games. As the sub-game policies are never explicitly stored, CFR-D then re-solves a sub-game on-demand using the CFR-D gadget.

This decomposition was initially motivated by the need to overcome memory limitations, essentially trading memory for time. The main ideas of the algorithm (decomposition and safe re-solving) now form the core building blocks of many modern search algorithms.

## 20.4   Sub-game Refinement



Figure 20.5: Sub-game refinement framework. (i) the strategy for the game is pre-computed using coarse-grained abstraction (ii) during the play, once we reach a node defining a sufficiently small subgame, we refine the strategy for that sub-game (iii) this together with the original strategy for the trunk creates a combined strategy. The point is to produce improved combined strategy

We now present another contribution of this thesis the — the subgame refinement framework and max margin re-solving. For a long time, state of the art methods simply pre-computed optimal solution for a simplified (abstracted) version of the game (Section 18.1). When the agent played the game, it simply acted based on this pre-computed policy. But as the game was approaching its end, the remaining sub-problem to reason about become either tractable or one could make a finer-grained abstraction of the sub-game.

An appealing idea is then to compute a new policy for the sub-game, refining the original strategy (Figure 20.5). One can than use the CFR-D gadget to guarantee the soundness of such approach. But even if one could do better in the finer-grained abstraction, CFR-D gadget has no strong reason to do so — it is only guaranteed to recover the original solution. That is, the resulting policy in the finer-grained abstraction is guaranteed to be no worse than the original policy based on the coarse abstraction.

## 20.4.1 Sub-game Margin

While the CFR-D gadget recovers the original policy, max margin aims to improve upon it and the resulting improvement is in a sense optimal (Moravcik et al., 2016). The safe re-solving constraints of CFR-D guaranteed that the opponent's counterfactual values of the resulting policy are no greater that counterfactual values of the original policy. To be able to improve the values, we first introduce the sub-game margin (Definition 25), corresponding to a slack between the values of the original policy and the values[2] of the new strategy. )

**Definition 25** (Sub-game Margin). *(Moravcik et al., 2016, Definition 2) Let $\pi_1, \pi_1'$ be a pair of strategies for sub-game $s_{pub}$. A sub-game margin $SM_1$ is then*

$$SM_1(\pi_1, \pi_1', s_{pub}) = \min_{s \in \mathcal{S}_2(s_{pub})} CBRV_2^{\pi_1}(s) - CBRV_2^{\pi_1'}(s)$$

The sub-game margin has several useful properties. The worst case performance is closely related to the value of the margin. If it is non-negative, the new strategy is guaranteed to be no more exploitable than original one. Furthermore, given that the opponent's best response reaches the sub-game with non-zero probability, the performance against this best response is improved and this improvement is at least proportional to the sub-game margin (and may be greater) (Theorem 29).

**Theorem 29.** *(Moravcik et al., 2016, Theorem 1) Given a strategy $\pi_1$, a sub-game $s_{pub}$ and a refined sub-game strategy $\pi_1^{s_{pub}}$, let $\pi_1' = \pi_1[s_{pub} \leftarrow \pi_1^{s_{pub}}]$ be a combined strategy of $\pi_1$ and $\pi_1^{s_{pub}}$. Let the sub-game margin $SM(\pi_1, \pi_1', s_{pub})$ be non-negative. Then $BRV(\pi_1') - BRV(\pi_1) \geq 0$. Furthermore, if there is a best response strategy $\pi_2^* \in \mathbb{BR}(\pi_1')$ such that $P^{(\pi_1', \pi_2^*)}(s) > 0$ for some $s \in \mathcal{S}_2(s_{pub})$, then $BRV(\pi_1') - BRV(\pi_1) \geq P_{-2}^{\pi_1'}(s) SM(\pi_1, \pi_1', s_{pub})$.*

*Proof.* See the Appendix of (Moravcik et al., 2016). □

Theorem 29 has critical consequences for the sub-game refinement. Since we refine the strategy once we reach the sub-game, we are either facing opponent's best response that reaches the $s_{pub}$, or they made a mistake earlier in the game. Furthermore, the probability of reaching a sub-game is proportional to $P^\pi$.

---

[2]Counterfactual best response values (Section 17.5.4).

## 20.4.2 Optimization Formulation

It is relatively easy to modify the CFR-D optimization formulation to maximize the sub-game margin by introducing a slack variable. The high level idea is that the CFR-D gadget construction simulates the "no greater than" constraints $u_s \leq V_{-i}(s_{pub}, s)$ (Section 20.2.1). Max margin then adds a slack/margin to the constraints that is then maximized $u_s + m \leq V_{-i}(s_{pub}, s)$ (Equation 20.2).

$$\max_{u,v,m} m$$
$$Fy = f, E^\top u - Ay \geq \mathbf{0}, y \geq \mathbf{0} \qquad\qquad (20.2)$$
$$u_s + m \leq V_{-i}(s_{pub}, s) \, \forall s \in \mathcal{S}_{-i}(s_{pub})$$

## 20.4.3 Gadget Game Construction



Figure 20.6: Comparison of the CFRD and max-margin gadgets. a) CFRD gadget: the opponent at each root infostate $s \in \mathcal{S}_{-i}(s_{pub})$ chooses to either (T)erminate, receiving the bound value $V_{-i}(s_{pub}, s)$, or to (F)ollow. b) Max-margin gadget: the opponent chooses an infostate $s \in \mathcal{S}_{-i}(s_{pub})$ to start with, and the corresponding terminal values are shifted so that the values under the original policy are all zero.

Similarly to the gadget-game construction of the CFR-D constraints (Section 20.2.2), we can create a gadget game corresponding to the max margin optimization. We create the gadget game by making two modifications to the original sub-game. i) we shift the opponent's utilities using the $CBRV_{-i}$, initializing all the state values to zero and ii) we add opponent's node followed by chance nodes at the top of the sub-game, allowing the opponent to pick any starting state.

We will distinguish the states, strategies, utilities, etc. for the gadget game by adding a tilde to corresponding notation. The following is a detailed description of the steps (see also Figure 20.6 that visualizes the constructed max margin gadget)

**Utilities Shift** To compare the changes in the performance of each of opponents' (root) information states, it is necessary to give them a common baseline. We use the original strategy $\pi_i^{s_{pub}}$ as the starting point. For every $s \in$

$\mathcal{S}_{-i}(s_{pub})$, we subtract the opponent's original counterfactual best response value from all the terminal histories $z$ reachable from $s$, setting the return to $\tilde{R}_{-i}(z) = R_{-i}(z) - CBRV_{-i}^{\pi_i}(s)$. We also update $\tilde{R}_i(z) = -\tilde{R}_{-i}(z)$ since we need the game to remain zero-sum. The reason behind this utilities shift is that under the original policy $\pi_i$, opponent's values are now zero for all of their root infostates: $CBRV_{-i}^{\pi_i}(\tilde{s}) = 0 \,\forall \tilde{s} \in \tilde{\mathcal{S}}_{-i}(s_{pu})$.

**Opponent's Node**  The opponent is permitted to choose the distribution of their private states at the start of the sub-game, while the player retains their distribution from the original strategy $\Delta(\mathcal{S}_i(s_{pub}))$. Since the opponent is aiming to maximize $\tilde{R}_{-i}$, they will always select the information state with the lowest margin. The minimax nature of the zero-sum game forces the player to find a strategy that maximizes this value. The construction adds a decision node $\tilde{d}$ for the opponent, where each action corresponds to choosing an information set $s \in \mathcal{S}_{-i}(s_{pub})$ to start with. Each action the leads to a new chance node $\tilde{c}_s$, where the chance player chooses the histories $h \in \mathcal{H}(s)$ based on $\Delta(\mathcal{S}_i(s_{pub}))$.

### 20.4.4   Experiments



Figure 20.7: Sub-game margins of the refined strategies. One big blind corresponds to 100 chips. The max-margin technique produces the optimal value. We see that the optimal value is much greater than the one produced by either re-solving or endgame solving (which produces even negative margins). The 95% confidence intervals for the results (after $10,000$ iterations) are: maxmargin $101.49 \pm 7.09$, re-solving $8.79 \pm 2.45$, unsafe solving $-518.5 \pm 49.19$.

We compare the unsafe re-solving, CFR-D re-solving and max-margin subgame refinement in sub-games of no-limit Texas hold'em poker. We use an improved version of the Nyx agent, the second strongest participant at the 2014 Annual Computer Poker Competition (heads-up no-limit Texas hold'em total bankroll) as the baseline strategy to be refined in the sub-games.

All the three techniques tested start with the same abstractions and trunk strategy. Following Ganzfried and Sandholm (2015), we begin the sub-games at

start of the last round (the river). While we use card abstraction to compute the original (trunk) strategy, specifically (Schmid et al., 2015) and (Johanson et al., 2013), the fine-grained abstraction for the endgame is using no card abstraction. We use the same actions in the refined sub-game as in the original strategy.

We refine only the sub-games that (after creating the fine-grained abstraction) are smaller than $1,000$ betting sequences — this is simply to speed up the experiments. The original agent strategy is used for both players in the trunk of the game. Once the self-play reaches the sub-game, we refine the strategy of the first player using each of the three techniques. We ran $10,000$ iterations of the CFR+ algorithm in the corresponding gadget games and each technique was used to refine around $2,000$ sub-games. Figure 20.7 then visualizes the average margins for the evaluated techniques.

**Unsafe Re-solving** The largely negative margin values suggest that the produced strategy may indeed be much more exploitable. **CFR-D Re-solving** The positive margin for re-solving shows that, although there's no explicit construction that forces the margin to be greater than zero, it does increase in practice. Notice, however, that the margin is far below the optimal level. **Max Margin Refinement** This technique produces a much larger sub-game margin than the previous techniques. The size of the margin suggests that the original strategy is potentially quite exploitable, and our technique can substantially decrease its exploitability (Theorem 29).

# 20.5   Full Lookahead Continual Resolving

Safe re-solving produces a policy for a single (public) sub-game and thus corresponds to a single step of search (Figure 20.1a), producing a strong globally consistent strategy. We now use safe re-solving as a building block in the continual re-solving algorithm.

As the name suggests, the algorithm continually performs re-solving step for the states visited during gameplay (Figure 20.1b). Continual re-solving thus needs to use and update the invariants required by the individual re-solving step. These invariants are updated on the game-play trajectory as the full lookahead rooted in the prior state includes the current game state and its required values.

## 20.5.1   Updating Invariants

Recall that the re-solving step requires the following invariants i) players' reach probabilities and ii) opponent's counterfactual values. First observe that for the very first state of the game, no re-solving is necessary and we can simply solve the initial game and store the result (Figure 20.8a). For a state visited during the game, we build a sub-game rooted in the current public state and use the corresponding invariants for the re-solve. The invariants are retrieved from the previous lookahead tree as it must include the currently visited state. This process is then repeated until the game is finished (Figure 20.8bc).

Figure 20.8: Continual Re-solving: a) Solve is run for the first state of the game. b) Re-solving step is performed for the next state the player is to act. To run re-solving, we construct the re-solving game using the opponent's value and player's reach probabilities. As the lookahead in the previous step included this state, the previous computation includes both of these quantities. c) We keep running re-solving for all the next nodes the player is to act in the same way.

### 20.5.2  Resulting Policy

As we need to reason about all the states in a public state, the algorithm produces a policy for all such states rather than the current infostate (i.e. strategy for all the poker hands we could be holding). Although policies for the states other than the current one are not directly useful for action selection, they are need to make a sound re-solve search. They can also be leveraged by variance reduction techniques (e.g. AIVAT, Section 21.6).

As the re-solving step uses the gadget game, any game solving method can be used. The common choice is to use a member of the CFR family, and continual re-solving has also been analysed with its Monte Carlo variants (Sustr et al., 2018).

### 20.5.3  Experiments

While the full lookahead is clearly impractical for large games, one can still evaluate the performance of the continual re-solving algorithm on small game variants. Figure 20.9 reports exploitability on Leduc poker and graph chase game (Glasses). As continual re-solving is a stateless online search algorithm (Section 7.3.1), we can compute the exploitability of the tabularized policy (Section  7.6).

## 20.6  Limited Lookahead and Value Functions

We are now ready to introduce a dept-limited solving. Just like in perfect information games, this allows us to reason over only a limited number of steps forward and evaluate the sub-games at the end of our lookahead using a value function. Fortunately, we already know how generalized value functions look for imperfect information games (Section 15.4). We just need to properly use them in the limited lookahead paradigm.

**Algorithm 7** Continual Re-solving with Full Lookahead

1: **function** PLAY($s \in \mathcal{S}_i$)
2:      $s_{pub} \leftarrow \mathcal{S}_{pub}(s)$
3:                        ▷ Uses the average strategy from the last search tree.
4:      $\Delta(\mathcal{S}_i(s_{pub})) \leftarrow GetReachProbabilities(last\_search\_tree, s_{pub})$
5:         ▷ Uses the averaged counterfactual values from the last search tree.
6:      $V_{-i}^{bound}(s_{pub}) \leftarrow GetCounterfactualValues(last\_search\_tree, s_{pub}, -i)$
7:      $search\_tree \leftarrow$ BUILDFULLLOOKAHEADTREE($s_{pub}$)
8:      SAFERESOLVE($search\_tree, \Delta(\mathcal{S}_i(s_{pub})), V_{-i}^{bound}(s_{pub}))$)
9:      $last\_search\_tree \leftarrow search\_tree$
10:      **return** $a \sim \pi_i(s)$
11:
12: **function** NEWGAME
13:      $last\_search\_tree \leftarrow$ BUILDFULLLOOKAHEADTREE($s_{pub}^{initial}$)
14:      SOLVE($last\_search\_tree$)



(a) Leduc poker.                               (b) Glasses.

Figure 20.9: Continual re-solving with full-lookahead tree and increasing number of CFR iterations for the search / re-solve.

The separation to lookahead tree and sub-games is possible thanks to the notion of public tree and we will show a simple modification of the public tree CFR that leverages value functions for the sub-games.

## 20.6.1    CFR and Value Functions

CFR is particularly suited to be modified to use value functions. This is because the CFR Theorem (Theorem 23) bounds the overall regret by the sum of individual counterfactual regrets. Furthermore, recall that the CFR algorithm (Algorithm 6) essentially sends reach probabilities down the tree and sends values up the tree. Consider what happens if we use a regret minimizer for all states in the lookahead tree, and use a zero-regret policy for all the states of the sub-games (Figure 20.10). The CFR Theorem then guarantees convergence of this approach.

As a zero-regret policy in a sub-game corresponds to a counterfactual best response, both players are best-responding to each other and thus playing optimally. During each CFR iteration, we can thus simply use an optimal policy in all the sub-games. Furthermore, the upward pass of CFR only requires values of

Figure 20.10: CFR with limited lookahead and value functions. During each iteration, regret minimizer is used to in the lookahead tree and zero-regret policies are played in the sub-games.

those sub-games. To run CFR in the lookahead tree, each iteration never explicitly requires one to compute the optimal (zero-regret) policies in the sub-games, it only requires the optimal values corresponding to sub-game — value functions.

Algorithm 8 then shows a minor modification to the ComputeValues Function (Algorithm 6, Line 10) that results in version of CFR that uses value functions. The reason why this is a particularly simple change is that our initial implementation of CFR was already operating on public states.

An important detail is that the counterfactual values returned by the value function need to be counterfactual best response values (Section 17.5.4). This is because value function returns values under an optimal policy, but we need to further restrict the corresponding optimal policy — making sure the per-iteration regrets in the sub-game are all zero. Fortunately, producing zero-regret optimal policies is not at all difficult and CFR algorithm already produces such policies.

---

**Algorithm 8** Limited Lookahead Tree CFR

---

1: **function** $\textsc{ComputeValues}(s_{pub} \in \mathcal{S}_{pub}, d_1 \in \Delta(\mathcal{S}_1(s_{pub})), d_2 \in \Delta(\mathcal{S}_2(s_{pub})))$
2:     **if** $\textsc{EndOfLookahead}(s_{pub})$ **then**
3:         $v_{i,c}(s) \leftarrow \textsc{ValueFunction}(s_{pub}, d_1, d_2)$
4:         **return**

---

## 20.7 Continual Re-solving with Value Functions

It is time to combine all of the building blocks introduced up to this point, culminating in a practical sound search algorithm for imperfect information games. We simply modify continual re-solving so that the individual re-solving steps use limited lookahead and generalized value functions (Figure 20.11). The combination of continual re-solving and value functions then results in the critical bound on the resulting policy (Theorem 30).

**Theorem 30.** *(Moravčík et al., 2017, Theorem 1) If the values returned by the value function used when the depth limit is reached have error less than $\epsilon$, and*

*T iterations of CFR are used to re-solve, then the exploitability of continual re-solving with value functions is less than $k_1\epsilon + k_2\sqrt{T}$, where $k_1$ and $k_2$ are game-specific constants.*

*Proof.* See the Appendix of Moravčík et al. (2017). □



Figure 20.11: Continual Resolving: a) Limited lookahead solve is run for the first state of the game. b) Re-solving step is performed for the next state the player is to act. To run re-solving, we construct the re-solving game using the opponent's value and player's reach probabilities. If the lookahead in the previous step included this state, the previous computation includes both of these quantities. c) We keep running re-solving for all the next nodes the player is to act in the same way.

### 20.7.1   Value Functions

Training generalized value functions is an open problem. Deep learning methods are a natural choice to learn and represent the generalized value functions, but how should we train these value functions? What training data and loss should one use? The ultimate goal is to learn value functions so that the search results in a low exploitable policy. But the exploitability is a complex function of the interaction of search and the value function. As this question is far from answered in perfect information games, it is hardly a surprise that there is little known at this point about the imperfect information case.

### 20.7.2   Experiments

We again use Leduc poker and Glasses to evaluate our search algorithm, and use tabularization for the exploitability computation. To see the effect of inexact value functions (i.e. functions producing value from an $\epsilon$-optimal policy with non-zero counterfactual regrets), we use the CFR algorithm with a varying number of iterations to serve as the value functions. This is important since in practice, we will have access only to an approximate value functions.

Figure 20.12 then shows the exploitability with a varying number of CFR iterations in search and in the value functions. We can see that as we increase the number of CFR iterations within the search tree, the exploitability gets lower

until it hits a limit of the imperfect value function. Improving the quality of the value function (by increasing number of CFR iterations for the value function) then allows the search to converge closer to an optimal solution. This is in line with the presented theorem, where the resulting exploitability is bound by the number of search iterations and the quality of the value function (Theorem 30).



(a) Leduc poker          (b) Graph chase game

Figure 20.12: Continual re-solving with limited lookahead tree. Lookahead tree is the smallest possible — only one step forward. For the value function, we use CFR to compute the target values with varying number of iterations (10, 20 and 40).

## 20.8    Limitations

Currently, there are some inherent limitations of the presented sound search methods. The re-solving steps as well as the value functions reason about all the states within a public state $\mathcal{S}_i(s_{pub})$. While necessary for sound re-solve methods, there are games where number of such states is intractably large. The methods also require explicit model of the environment (i.e. game tree) to run the search. Finally, as the model is discrete it can not capture continuous action space.

# 21. DeepStack

DeepStack — the final contribution of the thesis — was the first to introduce the combination of sound search and value functions in imperfect information game of poker (Moravčík et al., 2017). The combination of continual resolving and neural networks as value functions led to a leap improvement over the prior methods. First, unlike the abstraction based techniques, DeepStack was unexploitable by the local best response (Section 18.1.2). Second, DeepStack became the first program to beat professional human players in no-limit Texas hold'em poker.

## 21.1 Search / Resolving Step



Figure 21.1: Lookahed used on different streets, see also Table 21.1

DeepStack used continual resolving with limited lookaheads and value functions (Section 20.7). For the re-solving construction, it used a modified variant of the CFR-D gadget (Section 20.2.2) and fixed sparse lookahead tree (details in Section 21.2). The lookahead tree was built all the way to the next round (for pre-flop and flop infostates) or until the end of the game (for turn and river). This guaranteed that the lookahead included all the possible next states so that we could correctly update the invariants during the continual resolving process.

DeepStack then used a hybrid of CFR ad CFR+ for the policy computation — combining regret matching plus, simultaneous updates and uniform weighting of the average policy. Furthermore, it skipped the first half of iterations when averaging the policy and counterfactual values. The motivation is that the early iterations are often relatively poor (especially in the case of no-limit poker). Skipping iterations has been proven to be sound and is a commonly used trick. For more details on the search parameters, see Table 21.1 and Figure 21.1.

**Preflop** Preflop used $2,000$ iterations, skipped $1,000$ of them and the lookahead was built until the beginning of flop. During the initial skipped iterations, the auxiliary preflop network (at the end of the preflop) was used to speed up the computation. This allows to use a single value function call for each state at the end of the preflop, rather than enumerating all the possible (isomorphic) $22,100$ flops. For the final iterations that actually accumulate the the average strategy and counterfactual values, we did the expensive enumeration and flop evaluations as we need the counterfactual values for continual resolving.

**Flop** Preflop used $1,000$ iterations, skipped 500 of them and the lookahead was built until the beginning of turn.

**Turn** Preflop used $1,000$ iterations, skipped 500 of them and the lookahead was built all the way until the end of the game. To speed up the computation, we used a bucketed abstraction for the river sub-games within the lookahead tree (the same bucketing as used for neural net representation described in Section 21.3.2).

**River** Preflop used $1,000$ iterations, skipped 500 of them and the lookahead was built all the way until the end of the game. No abstraction other than the sparse lookaheads were used as these river sub-games can be resolved relatively quickly.

| Round | CFR Iterations | Skip Iterations | Lookahead | Value Functions |
|---|---|---|---|---|
| Pre-flop | 1,000 | 500 | Until Flop | Flop/Aux |
| Flop | 1,000 | 500 | Until Turn | Turn |
| Turn | 1,000 | 500 | Full | - |
| River | 2,000 | 1,000 | Full | - |

Table 21.1: Detailed parameters of the re-solving step of DeepStack (see also Figure 21.1).

## 21.2   Action Abstraction in Lookahead



Figure 21.2: (a) Resolving a state included in the abstraction. We simply use the invariants for that state. (b) Resolving a state after the opponent made an action not included in the sparse lookahead. Due to structure of the game and the lookaheads used by DeepStack, we can use invariants from the previous state.

DeepStack used sparse lookahead trees, where it considered only a subset of all the possible actions (there are almost $20,000$ actions in each infostate). Sparse lookahead trees are now a common choice in online search agents for the no-limit

| Poker Round | First Actions | Second Actions | Remaining Actions |
|---|---|---|---|
| Pre-flop | F, C, ½P, P, A | F, C, ½P, P, 2P, A | F, C, P, A |
| Flop | F, C, ½P, P, A | F, C, P, A | F, C, P, A |
| Turn | F, C, ½P, P, A | F, C, P, A | F, C, P, A |
| River | F, C, ½P, P, P, A | F, C, ½P,P, P, A | F, C, P, A |

Table 21.2: Sparse lookahead actions. (F)old, (C)all/(C)heck, (P)ot bet and fractional (P)ot bets, (A)ll-in.

poker (Brown and Sandholm, 2018; Brown et al., 2020; Zarick et al., 2020). For the full list of actions included in DeepStack's search see Table 21.2.

There is an important distinction to the abstraction techniques described in Section 18. DeepStack never needs to translate the current (real) state into an abstraction, regardless of the actions taken by the opponent. The resolving step always starts in the exact state of the game, and DeepStack thus perfectly understands the current situation and the pot-ods. This is confirmed by our LBR analysis in Section 21.5.

But to start the safe resolving step for $s_{pub}$, we still require the counterfactual values of the opponent $V_{-i}^{bound}(s_{pub})$ and agent's reach distribution $\Delta(\mathcal{S}_i(s_{pub}))$ (Section 20.2). DeepStack leveraged a particular structure of the game to correctly compute these invariants even if the current state $s_{pub}$ was not included in the sparse search tree.

**Counterfactual Values** First, consider the counterfactual values. As there is always a single player to act and all the actions are publicly observable, we can assign a unique player to each $s_{pub}$ (agent/opponent/chance). Under any best responding policy, the following then holds for any $s_{pub}$ of the opponent $\forall s \in \mathcal{S}_{-i}(s_{pub}) \forall a \in \mathcal{A}(s) : v_{-i,c}^{\pi}(s) \geq q_{-i,c}^{\pi}(s,a)$. Furthermore, there is a single state $sa$ after an action $a$ in a state $s$. We thus have $v_{-i,c}^{\pi}(s) \geq v_{-i,c}^{\pi}(sa)$. This allows us to use counterfactual values of the state $s$ as the constraint values for $sa$ — that is, we only need the values for the state prior to the action of the opponent. DeepStack's lookahead tree then guarantees that this is the only off-tree case. This is because the lookahead included all the possible next chance states and DeepStack will never take an actions not included in its abstractions.

**Reach Distribution** Agent's distribution over their infostates $\Delta(\mathcal{S}_i(s_{pub}))$ is not a function of the opponent's policy and we can thus compute the reach distribution for the current state $s_{pub}$ after we observe the state (i.e. after the opponent made an action).

### 21.2.1 Action Abstraction Analysis

The choice of sparse lookahead tree abstraction thus only effects i) the actions the agent can make and ii) the quality of the counterfactual values to be used in the next re-solving step. While we know that it is sound to use the values of the parent state $s$ for resolving a state $sa$, we can only do so because we assumed the

| Betting | Size | $L_1$ | $L_2$ | $L_\infty$ |
|---|---|---|---|---|
| F, C, Min, ¹/₄P, ¹/₂P, ³/₄P, P, 2P, 3P, 10P, A | 555k | 18.06 | 0.891 | 0.2724 |
| F, C, P, A | 61k | 25.51 | 1.272 | 0.3372 |
| F, C, 2P, A | 48k | 64.79 | 2.672 | 0.3445 |
| F, C, ¹/₂P, A | 100k | 58.24 | 3.426 | 0.7376 |
| F, C, ¹/₂P, P, A | 126k | 41.42 | 1.541 | 0.2955 |
| F, C, P, 2P, A | 204k | 27.69 | 1.390 | 0.2543 |
| F, C, ¹/₂P, P, 2P, A | 360k | 20.96 | 1.059 | 0.2653 |

Table 21.3: Errors in counterfactual values when solving a river game with different action abstraction. The counterfactual values computed with $1,000$ and averaged over 100 random river situations. As the "ground truth" targets, we used the largest set of actions (F, C, Min, ¹/₄P, ¹/₂P, ³/₄P, P, 2P, 3P, 10P, A) with $4,000$ iterations.

values for $s$ are based on a best responding policy. As we limit the actions in the lookahead, the resulting values might be considerably different.

The action set used by DeepStack was constructed using a careful analysis, where experimented with different abstraction sizes and measured the quality of the resulting counterfactual values (Table 21.3). Additional experiments proved that it is critical to include fold, call, all-in and a bet around the pot size — and the same choice is then made by other search agents (Section 21.7).

## 21.3 Value Function

DeepStack used deep fully connected neural networks with seven hidden layers to represent the value functions. Each layer contained 500 neurons and used the parametric rectified linear units (He et al., 2015). The output of the final linear layer then predicts individual counterfactual values and thus contains as many neurons as there are infostates. Since the game is zero-sum, it must be the case that the reach-weighted sum of counterfactual values is in balance (Equation 21.1). To make sure this is the case, we included additional custom layer that redistributes the "excess", ensuring the network produces only zero-sum values.

$$\sum_{s \in \mathcal{S}_i(s_{pub})} \Delta(\mathcal{S}_i(s_{pub})(s)v_{i,c}(s) + \sum_{s \in \mathcal{S}_{-i}(s_{pub})} \Delta(\mathcal{S}_{-i}(s_{pub})(s)v_{-i,c}(s) = 0 \qquad (21.1)$$

We trained separate networks for flop and turn that were used during preflop and flop search respectively (as the lookahed reached all the way to the next street). Furthermore, an auxiliary preflop network was used to speed up the computation on preflop.

### 21.3.1 Networks Training

All the networks were trained using the supervised training paradigm. The training data was generated by sampling random[1] sub-games (inputs) and solving them to obtain the counterfactual values (outputs/targets). These input-output pairs were then used as the supervised dataset. As the lookahead was built all the way to the next street (Section 21.1), the sampled sub-games always corresponded to the initial public states of the respective streets (with the exception of the auxiliary preflop network where the public states rather corresponded to the last states of preflop).

The solving used $1,000$ iterations of CFR+ and the FCPA action abstraction. The training was implemented using the Torch7 libraries (Collobert et al., 2011). The training loss was the average Huber loss (Huber, 1992) over the counterfactual values, minimized using the Adam stochastic gradient descent procedure (Kingma and Ba, 2014). We used batch size of $1,000$ and a learning rate $0.001$, which was decreased to $0.0001$ after the first 200 epochs. Networks were trained for approximately 350 epochs over two days on a single GPU, and the epoch with the lowest validation loss was then selected.

**Turn** The training data consisted of ten million randomly sampled turn sub-games that were then solved.

**Flop** We sampled one million flop sub-games and used limited lookahead solving (Section 20.6) with the already trained turn network as the value function — thus bootstrapping the networks.

**Aux Preflop** We sampled ten million preflop situations and the target values were obtained by enumerating all $22,100$ possible flops and averaging the counterfactual values from the flop network's output (another level of bootstrapping). This averaging is equivalent to the limited lookahead solving used for the flop network, as there are only chance states in between the last preflop state and the flop network.

### 21.3.2 Networks Features

To simplify the generalization task of the networks, we map the distribution over the individual poker hands (combinations of public and private cards) into distribution over buckets/clusters. This essentially compresses the state space of $\binom{52}{7}$ card combinations down to the number of buckets used. For both the turn and flop networks, we used $1,000$ clusters generated using k-means clustering with earth mover's distance over hand-strength features (Ganzfried and Sandholm, 2014; Johanson et al., 2013). The auxiliary preflop network used no bucketing as there are only 169 isomorphic hands.

---

[1]See supplementary material of DeepStack for details on the distribution.

## 21.4 Human Evaluation

In collaboration with the International Federation of Poker (IFP, now the International Federation of Match Poker IFMP) (IFMP, 2021), we recruited 33 players from 17 countries. Each player was expected to play $3,000$ games in between November 7th and December 12th, 2016. Cash incentives were given to the top three performers ($\$5,000$, $\$2,500$ and $\$1,250$ CAD).

Poker is inherently game of high variance and even if one agent is substantially stronger than the opponent, one might have to play hundreds of thousands of games to get statistically significant result. We thus evaluate the performance using AIVAT (Section 21.6), where the counterfactual values produced by the continual resolving served as the value estimates required by AIVAT — providing an excellent estimates and resulting in impressive 85% reduction of standard deviation.

A total of $44,852$ games were finished by the players, with 11 players completing the full $3,000$ games. DeepStack won $492 \pm 220$ mbb/g when the raw data was used, and when AIVAT correction terms were used the performance was estimated at $486$ mbb/g $\pm 40$ mbb/g (note the significantly smaller confidence interval). AIVAT also allows us to derive statistically significant individual results for all but one of the players who finished the required $3,000$ games. DeepStack is beating all of such player, with 10 out of 11 results being significant and only for the best performing player is the result not significant. The detailed results can be found in Table 21.4 and Figure 21.3



Figure 21.3: Win rate of DeepStack against the individual human players: AIVAT value with 95% confidence interval and number of hands per player. Thanks to AIVAT variance reduction, the result is significant for all but one individual player who finished the required $3,000$ hands. See also Table 21.4 for more analysis and aggregate results.

| Player | Games | Win Rate (mbb/g) |
|---|---|---|
| Martin Sturc | 3000 | $70 \pm 119$ |
| Stanislav Voloshin | 3000 | $126 \pm 103$ |
| Prakshat Shrimankar | 3000 | $139 \pm 97$ |
| Ivan Shabalin | 3000 | $170 \pm 99$ |
| Lucas Schaumann | 3000 | $207 \pm 87$ |
| Phil Laak | 3000 | $212 \pm 143$ |
| Kaishi Sun | 3000 | $363 \pm 116$ |
| Dmitry Lesnoy | 3000 | $411 \pm 138$ |
| Antonio Parlavecchio | 3000 | $618 \pm 212$ |
| Muskan Sethi | 3000 | $1009 \pm 184$ |
| Pol Dmit | 3000 | $1008 \pm 156$ |
| Tsuneaki Takeda | 1901 | $627 \pm 231$ |
| Youwei Qin | 1759 | $1306 \pm 331$ |
| Fintan Gavin | 1555 | $635 \pm 278$ |
| Giedrius Talacka | 1514 | $1063 \pm 338$ |
| Juergen Bachmann | 1088 | $527 \pm 198$ |
| Sergey Indenok | 852 | $881 \pm 371$ |
| Sebastian Schwab | 516 | $1086 \pm 598$ |
| Dara O'Kearney | 456 | $78 \pm 250$ |
| Roman Shaposhnikov | 330 | $131 \pm 305$ |
| Shai Zurr | 330 | $499 \pm 360$ |
| Luca Moschitta | 328 | $444 \pm 580$ |
| Stas Tishekvich | 295 | $-45 \pm 433$ |
| Eyal Eshkar | 191 | $18 \pm 608$ |
| Jefri Islam | 176 | $997 \pm 700$ |
| Fan Sun | 122 | $531 \pm 774$ |
| Igor Naumenko | 102 | $-137 \pm 638$ |
| Silvio Pizzarello | 90 | $1500 \pm 2100$ |
| Gaia Freire | 76 | $369 \pm 136$ |
| Alexander Bös | 74 | $487 \pm 756$ |
| Victor Santos | 58 | $475 \pm 462$ |
| Mike Phan | 32 | $-1019 \pm 2352$ |
| Juan Manuel Pastor | 7 | $2744 \pm 3521$ |

Table 21.4: DeepStack's win rate against professional poker players, with 95% confidence interval.

### 21.4.1 Thinking Time

To minimize the thinking time, both search and network evaluation were ran on GPU (single NVIDIA GeForce GTX 1080) and we batched the network evaluations. To further speed up the play, we cached the resolving result for every observed preflop situation. When the same information state was reached again, we simply retrieve the result from the cache. Table 21.5 then reports thinking times for DeepStack and humans. Note that DeepStack acted considerably faster than our human players in all rounds, and that the pre-flop round was by far the fastest due to the cache.

| | Humans | | DeepStack | |
| Round | Median | Mean | Median | Mean |
| --- | --- | --- | --- | --- |
| Pre-flop | 10.3 | 16.2 | 0.04 | 0.2 |
| Flop | 9.1 | 14.6 | 5.9 | 5.9 |
| Turn | 8.0 | 14.0 | 5.4 | 5.5 |
| River | 9.5 | 16.2 | 2.2 | 2.1 |
| Per Action | 9.6 | 15.4 | 2.3 | 3.0 |
| Per Hand | 22.0 | 37.4 | 5.7 | 7.2 |

Table 21.5: Thinking time [s] of DeepStack and humans.

## 21.5 Local Best Response Evaluation

Not only DeepStack was the first to beat professional human in no-limit Texas hold'em poker. Unlike the previous approaches based on the abstraction framework (Section 18.1), suggesting that it's substantially harder to exploit and closer to optimal policy. Furthermore, we can evaluate how the choice of the actions in the sparse lookahead affect the LBR performance. Table 21.7 reports the LBR results with three different actions abstractions. Regardless of the abstraction in place, LBR is unable to exploit the agent.

## 21.6 AIVAT Variance Reduction Technique

No-limit poker is inherently noisy game, where one needs to collect a large number of games to draw statistical conclusions. Consider the first human-machine event in no-limit Texas hold'em poker, where the Claudico agent faced a team of four professional poker players (Ganzfried, 2016). The match involved $80,000$ hands of poker, where each human played dozens of hours over the span of seven days. Despite this significant investment of human time and Claudico losing by over $90mbb$, the result was still not statistically conclusive (it was on the very edge of the statistical test). It is wort mentioning that Claudico was an instance of abstraction based agent (Section 18.1) combined with unsafe resolving method (Section 20.1) used for sub-game refinement.

| | | | | |
|---|---|---|---|---|
| LBR Pre-flop Actions | F, C | C | C | C |
| LBR Flop Actions | F, C | C | C | 56bets |
| LBR Turn Actions | F, C | F, C, P, A | 56bets | F, C |
| LBR River Actions | F, C | F, C, P, A | 56bets | F, C |
| Hyperborean (2014) | $721 \pm 56$ | $3852 \pm 141$ | $4675 \pm 152$ | $983 \pm 95$ |
| Slumbot (2016) | $522 \pm 50$ | $4020 \pm 115$ | $3763 \pm 104$ | $1227 \pm 79$ |
| Act1 (2016) | $407 \pm 47$ | $2597 \pm 140$ | $3302 \pm 122$ | $847 \pm 78$ |
| Full Cards [100 BB] | $-424 \pm 37$ | $-536 \pm 87$ | $2403 \pm 87$ | $1008 \pm 68$ |
| DeepStack | $-428 \pm 87$ | $-383 \pm 219$ | $-775 \pm 255$ | $-602 \pm 214$ |

Table 21.6: Local best response performance against strong abstraction-based poker agents in no-limit Texas hold'em poker [mbb/g]. The list of actions considered by the technique varies in different poker rounds, and we report four different configurations. Actions are (F)old, (C)all (P)ot bet and (A)ll-in. For the full list of 56 bets see (Lisy and Bowling, 2017). Note that a policy that simply always folds each hand is exploitable by 750*mbb*.

| First level actions | LBR performance |
|---|---|
| F, C, P, A | $-479 \pm 216$ |
| Default | $-383 \pm 219$ |
| F, C, ⅓P, P, ½P, 2P, A | $-406 \pm 218$ |

Table 21.7: LBR against DeepStack with different action sets in DeepStack's lookahead tree.

Another large scale human match took place from January 11 to January 31, 2017. Libratus (Section 21.7.1) faced four top-class human poker players (Jason Les, Dong Kim, Daniel McAulay and Jimmy Chou). While this time the results allowed for statisticaly significant conclusion, the match involved a total of $120,000$ hands ($50\%$ increase compared to Claudico). And since this time Libratus was (mostly) a search based agent, it bested its human opponents.

## 21.6.1  Previous Variance Reduction Methods

Duplicate (or mirror) hands is a simple idea aimed to decrease the role of luck introduced by the chance events (Davidson, 2014). Mirror hands methods replicate the chance events for both positions of the players. Namely in poker, the agents get to play both positions (small and big blind seats) with the identical chance deals (both private and public cards). The motivation is that if one player gets "lucky" by favourable card deal, the variance will be reduced as the agent also gets to play from the less favourable position (note that this method also halves the number of datapoints as the paired outcome then forms a single measurement). This is of course problematic when human play is involved, as we need to reset

the memory of the agent to make sure they can not predict the chance events. Typical workaround is then to play against human pairs, where the paired humans are dealt the duplicate hands while playing in a separate room (room one: computer vs player one, room two: computer vs player two) (Hedberg, 2007). This pairing method was indeed used during both Claudico and Libratus events. While simple, this method provides only a modest improvement.

Second, the technique of importance sampling over imaginary observations (Bowling et al., 2008) uses the knowledge of player's policy to replace the single terminal state return $R_i(z)$ with expectation over the compatible terminal states in $\mathcal{S}_{pub}(z)$ and uses the importance sampling correction. In poker terms, this replaces the single hand outcome with the expected value over all the other hands the agent could be dealt and their respective reach probabilities to that terminal public state (often referred to as the range evaluation in poker literature).

Third, techniques like MIVAT (White and Bowling, 2009) use the method of control variates (Section 17.8.1) to reduce the variance caused by chance events. Each chance action is paired with a value estimate (baseline), resulting in correction terms for each chance event on the trajectory (just like in VR-MCCFR in Section 17.8). In poker, these correction terms essentially aim to negate the chance "luck" factor, e.g. when the agent hits good/bad card deal.

## 21.6.2   AIVAT



(a)          (b)          (c)          (d)

Figure 21.4: (a) Monte Carlo estimate uses the single sampled outcome $R_i(z)$. (b) Importance sampling over imaginary observations considers the outcomes of all the possible private states. (c) MIVAT uses a single outcome and control variate with baselines to form correction terms for each chance event on the trajectory (d) AIVAT uses imaginary observations in combination with correction terms applied for chance events as well as the the agent's decitions.

AIVAT (action-informed value assessment tool, (Burch et al., 2018)) combines the technique of imaginary observations and control variates. Unlike MIVAT, it uses the control variate based correction terms both for the chance actions and the agent actions. There are no correction terms for opponent's actions as we do not know their policy.

We compare the techniques in no-limit Texas hold'em poker, using a relatively small abstraction agent trained with MCCFR (Section 17.7) for the underlying policy. The resulting counterfactual values are then used as the baselines for the correction terms. The data are generated from self-play of the agent, using

| Estimator | $\bar{v}_x$ | $SD(v_x)$ |
|---|---|---|
| Monte Carlo (chips) | 0.03871 | 25.962 |
| MIVAT | 0.02038 | 21.293 |
| MIVAT+Imaginary | 0.02596 | 16.073 |
| AIVAT | 0.00186 | 8.095 |

Table 21.8: Comparison of variance reduction techniques in no-limit Texas hold'em poker. Data are computed from self-play of a small abstraction agent for a 100 big blind variant.

1 million games. Table 21.8 compares the variance of Monte Carlo estimator, MIVAT and AIVAT.

AIVAT leads to substantially lower variance compared to prior techniques, resulting in 68% decrease in standard deviation. Just like better value estimates lead to lower variance for VR-MCCFR, so do they in the case of AIVAT. While this is substantial improvement over the prior techniques, the same variance reduction technique in DeepStack proved even more effective. This suggests that valued produced by the continual re-solving are much closer to the true values despite the fact that DeepStack faced human players.

## 21.7 Other Search Agents

After DeepStack, there are now multiple agents that successfully employ search techniques in imperfect information games.

### 21.7.1 Libratus

Shortly after DeepStack, Libratus used search to beat top human professional players (Brown and Sandholm, 2018). Libratus has essentially two phases. When the sub-game is small enough to be reasoned about, Libratus runs continual resolving with full lookahead and sparse action abstraction (the authors refer to the continual resolving as "nested resolving" as they developed the technique independently in parallel (Brown and Sandholm, 2017)). In the earlier rounds though, Libratus did not use limited lookaheads combined with a value function. Rather than doing search, it falls back to the prior abstraction based techniques (Section 18), where the pre-computed abstraction policy is referred to as the "blue print" policy.

### 21.7.2 Modicum

Following Libratus, the authors presented Modicum (Brown et al., 2018b). Modicum replaces the "blue print" policy and uses limited lookahead search with value functions. Rather than learning the value function, a hand-crafted rollout-based evaluation function is used. Advantage of this evaluation function is that

it is particularly fast and easy to evaluate. On the other hand, it has very limited guarantees and one has to construct this evaluation function for each game manually

### 21.7.3  Supremus

Supremus is essentially a re-implementation of DeepStack with larger neural networks, larger search tree and more training data (Zarick et al., 2020). Supremus was then involved during the 2021 public match between Doug Paulk (considered to be the best heads-up no-limit Texas hold'em specialist in the world) and Daniel Negreanu (recognized as the best poker player of the decade in 2014 and having won over $42,000,000$ USD). Supremus was used by Doug Paulk to help him to formulate the strategies, and he ended up winning around $1,200,000$ USD (Hedberg, 2021). Thanks to the continual re-solving and learned value functions, state of the art agents in few years improved from being easily exploitable by simple techniques, to beating even the best human players.

### 21.7.4  DeepRole

DeepRole is an agent for the game "The Resistance: Avalon" — a multi-agent hidden role game. Its search algorithm combines CFR, search and value functions operating on public states, trained similarly to DeepStack (Serrino et al., 2019).

### 21.7.5  ReBeL

Rebel is a recent algorithm that can be though of as "simulating" CFR-D (Section 20.3) (Brown et al., 2020). While the algorithm does not have to use gadget during the re-solving, it comes at a price. It makes very strong assumptions about the re-solving policy, where it assumes an exact match between its training and evaluation phase. Namely, the re-solving policy in a sub-game is assumed to match the one that produced the values during the trunk computation. In other words, rather than using gadget, is uses the unsafe resolve (Section 20.1) and assumes the resulting policy to match the one that originally produced the sub-game values. This strong assumptions also limits the algorithm to use different lookaheads and thinking time during the actual play (a crucial property of strong search based agents in perfect information games), as it introduces a mismatch between the training and evaluation phase.

This algorithm thus have fundamentally different properties and assumptions to the agents built on the notion of continual re-solving, where the only assumption is quality of the value function

# Conclusion

Thanks to the combination of search and (learned) value functions, computers bested even the strongest of human players in Chess, Go, Backgammon, Arima and many other perfect information games. Generalization of the underlying concepts to imperfect information then allows for search methods in imperfect information, resulting in agents that outplay humans in poker — challenging game of imperfect information.

Imperfect information games allow for modeling much wider class of problems and interactions. As many real world problems do not allow for perfect information, the concepts introduced in this thesis allow for potentially interesting applications.

# Bibliography

Arjona-Medina, J. A., Gillhofer, M., Widrich, M., Unterthiner, T., Brandstetter, J., and Hochreiter, S. (2018). Rudder: Return decomposition for delayed rewards. *arXiv preprint arXiv:1806.07857*.

Astrom, K. J. (1965). Optimal control of markov processes with incomplete state information. *Journal of mathematical analysis and applications*, 10(1):174–205.

Aumann, R. J. (1961). Mixed and behavior strategies in infinite extensive games. Technical report, PRINCETON UNIV NJ.

Bailey, J. P. and Piliouras, G. (2018). Multiplicative weights update in zero-sum games. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 321–338.

Bard, N. D. (2016). *Online Agent Modelling in Human-Scale Problems*. PhD thesis, University of Alberta.

Bellman, R. (1957). A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684.

Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C., et al. (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*.

Bernstein, D. S., Givan, R., Immerman, N., and Zilberstein, S. (2002). The complexity of decentralized control of markov decision processes. *Mathematics of operations research*, 27(4):819–840.

Bertsimas, D. and Tsitsiklis, J. N. (1997). *Introduction to linear optimization*, volume 6. Athena Scientific Belmont, MA.

Blackwell, D. et al. (1956). An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics*, 6(1):1–8.

Bosansky, B., Kiekintveld, C., Lisy, V., and Pechoucek, M. (2014). An exact double-oracle algorithm for zero-sum extensive-form games with imperfect information. *Journal of Artificial Intelligence Research*, 51:829–866.

Bowling, M. (2003). *Multiagent learning in the presence of agents with limitations*. PhD thesis, CARNEGIE-MELLON UNIV PITTSBURGH PA SCHOOL OF COMPUTER SCIENCE.

Bowling, M., Burch, N., Johanson, M., and Tammelin, O. (2015). Heads-up limit hold'em poker is solved. *Science*, 347(6218):145–149.

Bowling, M., Johanson, M., Burch, N., and Szafron, D. (2008). Strategy evaluation in extensive games with importance sampling. In *Proceedings of the 25th international conference on Machine learning*, pages 72–79.

Boyd, S. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.

Brown, G. W. (1951). Iterative solution of games by fictitious play. *Activity analysis of production and allocation*, 13(1):374–376.

Brown, N., Bakhtin, A., Lerer, A., and Gong, Q. (2020). Combining deep reinforcement learning and search for imperfect-information games. *Advances in Neural Information Processing Systems*, 33.

Brown, N., Lerer, A., Gross, S., and Sandholm, T. (2018a). Deep counterfactual regret minimization. *arXiv preprint arXiv:1811.00164*.

Brown, N. and Sandholm, T. (2017). Safe and nested subgame solving for imperfect-information games. In *Advances in neural information processing systems*, pages 689–699.

Brown, N. and Sandholm, T. (2018). Superhuman ai for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424.

Brown, N., Sandholm, T., and Amos, B. (2018b). Depth-limited solving for imperfect-information games. In *Advances in Neural Information Processing Systems*, pages 7663–7674.

Browne, C. B., Powley, E., Whitehouse, D., Lucas, S. M., Cowling, P. I., Rohlfshagen, P., Tavener, S., Perez, D., Samothrakis, S., and Colton, S. (2012). A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43.

Burch, N. (2018). *Time and space: Why imperfect information games are hard.* PhD thesis, University of Alberta.

Burch, N., Johanson, M., and Bowling, M. (2014). Solving imperfect information games using decomposition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28.

Burch, N., Moravcik, M., and Schmid, M. (2019). Revisiting cfr+ and alternating updates. *Journal of Artificial Intelligence Research*, 64:429–443.

Burch, N., Schmid, M., Moravcik, M., Morill, D., and Bowling, M. (2018). Aivat: A new variance reduction technique for agent evaluation in imperfect information games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.

Campbell, M., Hoane Jr, A. J., and Hsu, F.-h. (2002). Deep blue. *Artificial intelligence*, 134(1-2):57–83.

Čermák, J., Bošansky, B., and Lisy, V. (2017). An algorithm for constructing and solving imperfect recall abstractions of large extensive-form games. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 936–942.

Chaudhuri, K., Freund, Y., and Hsu, D. (2009). A parameter-free hedging algorithm. *arXiv preprint arXiv:0903.2851*.

Chen, B. and Ankenman, J. (2006). *The mathematics of poker.* ConJelCo LLC.

Chen, X. and Deng, X. (2006). Settling the complexity of two-player nash equilibrium. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, pages 261–272. IEEE.

Chiang, C.-K., Yang, T., Lee, C.-J., Mahdavi, M., Lu, C.-J., Jin, R., and Zhu, S. (2012). Online optimization with gradual variations. In *Conference on Learning Theory*, pages 6–1. JMLR Workshop and Conference Proceedings.

Collobert, R., Kavukcuoglu, K., and Farabet, C. (2011). Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*.

Copeland, B. J. (2004). *The essential turing*. Clarendon Press.

Daskalakis, C., Deckelbaum, A., and Kim, A. (2011). Near-optimal no-regret algorithms for zero-sum games. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 235–254. SIAM.

Daskalakis, C., Goldberg, P. W., and Papadimitriou, C. H. (2009). The complexity of computing a nash equilibrium. *SIAM Journal on Computing*, 39(1):195–259.

Davidson, J. (2014). The baseline approach to agent evaluation. Master's thesis, University of Alberta.

Davis, T., Schmid, M., and Bowling, M. (2020). Low-variance and zero-variance baselines for extensive-form games. In *International Conference on Machine Learning*, pages 2392–2401. PMLR.

Farina, G., Kroer, C., and Sandholm, T. (2018). Online convex optimization for sequential decision processes and extensive-form games. *arXiv preprint arXiv:1809.03075*.

Farina, G., Kroer, C., and Sandholm, T. (2019a). Online convex optimization for sequential decision processes and extensive-form games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 1917–1925.

Farina, G., Kroer, C., and Sandholm, T. (2019b). Optimistic regret minimization for extensive-form games via dilated distance-generating functions. *arXiv preprint arXiv:1910.10906*.

Farina, G., Kroer, C., and Sandholm, T. (2020). Faster game solving via predictive blackwell approachability: Connecting regret matching and mirror descent. *arXiv preprint arXiv:2007.14358*.

Frank, I., Basin, D. A., and Matsubara, H. (1998). Finding optimal strategies for imperfect information games. In *AAAI/IAAI*, pages 500–507.

Freund, Y. and Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139.

Gale, D., Kuhn, H., and Tucker, A. (1951). Linear programming and the theory of games. In et al., T. K., editor, *Activity Analysis of Production and Allocation*, pages 317–329. Wiley: New York.

Ganzfried, S. (2016). Reflections on the first man vs. machine no-limit texas hold'em competition. *ACM SIGecom Exchanges*, 14(2):2–15.

Ganzfried, S. and Sandholm, T. (2014). Potential-aware imperfect-recall abstraction with earth mover's distance in imperfect-information games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28.

Ganzfried, S. and Sandholm, T. (2015). Endgame solving in large imperfect-information games. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 37–45. Citeseer.

Gibson, R., Lanctot, M., Burch, N., Szafron, D., and Bowling, M. (2012). Generalized sampling and variance in counterfactual regret minimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26.

Gilpin, A. and Sandholm, T. (2007). Lossless abstraction of imperfect information games. *Journal of the ACM (JACM)*, 54(5):25–es.

Greenwald, A., Li, J., and Sodomka, E. (2017). Solving for best responses and equilibria in extensive-form games with reinforcement learning methods. In *Rohit Parikh on Logic, Language and Society*, pages 185–226. Springer.

Gruslys, A., Lanctot, M., Munos, R., Timbers, F., Schmid, M., Perolat, J., Morrill, D., Zambaldi, V., Lespiau, J.-B., Schultz, J., et al. (2020). The advantage regret-matching actor-critic. *arXiv preprint arXiv:2008.12234*.

Hansen, E. A., Bernstein, D. S., and Zilberstein, S. (2004a). Dynamic programming for partially observable stochastic games. In *AAAI*, volume 4, pages 709–715.

Hansen, E. A., Bernstein, D. S., and Zilberstein, S. (2004b). Dynamic programming for partially observable stochastic games. In *AAAI*, volume 4, pages 709–715.

Hansen, K. A., Miltersen, P. B., and Sørensen, T. B. (2007). Finding equilibria in games of no chance. In *International Computing and Combinatorics Conference*, pages 274–284. Springer.

Harris, C. (1998). On the rate of convergence of continuous-time fictitious play. *Games and Economic Behavior*, 22(2):238–259.

Harroch, R. D. and Krieger, L. (2007). *Poker for dummies*. John Wiley & Sons.

Hart, S. and Mas-Colell, A. (2000). A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68(5):1127–1150.

Hawkin, J., Holte, R., and Szafron, D. (2011). Automated action abstraction of imperfect information extensive-form games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25.

He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034.

Hedberg, S. (2007). Man vs. machine poker challenge. `https://www.aaai.org/Pressroom/Releases/release-07-0612.pdf`. [Online; accessed 20-May-2021].

Hedberg, S. (2021). Team polk's bryan pellegrino talks about his ai research and how it helped formulate strategies to win \$1.2 million. `https://www.cardplayer.com/poker-news/25778-team-polk-s-bryan-pellegrino-talks-about-his-ai-research-\and-how-it-helped-formulate-strategies-to-win-1-2-million`. [Online; accessed 23-May-2021].

Heinrich, J., Lanctot, M., and Silver, D. (2015). Fictitious self-play in extensive-form games. In *International conference on machine learning*, pages 805–813. PMLR.

Herbster, M. and Warmuth, M. K. (1998). Tracking the best expert. *Machine learning*, 32(2):151–178.

Huber, P. J. (1992). Robust estimation of a location parameter. In *Breakthroughs in statistics*, pages 492–518. Springer.

IFMP (2021). The international federation of match poker. `https://matchpokerfed.org/`. [Online; accessed 06-May-2021].

Jakobsen, S. K., Sørensen, T. B., and Conitzer, V. (2016). Timeability of extensive-form games. In *Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science*, pages 191–199.

Johanson, M. (2013). Measuring the size of large no-limit poker games. *arXiv preprint arXiv:1302.7008*.

Johanson, M., Burch, N., Valenzano, R., and Bowling, M. (2013). Evaluating state-space abstractions in extensive-form games. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*, pages 271–278.

Johanson, M., Waugh, K., Bowling, M., and Zinkevich, M. (2011). Accelerating best response calculation in large extensive games. In *IJCAI*, volume 11, pages 258–265.

Johanson, M. B. (2016). *Robust strategies and counter-strategies: from superhuman to optimal play*. PhD thesis, University of Alberta.

Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. (1998). Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134.

Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Knuth, D. E. and Moore, R. W. (1975). An analysis of alpha-beta pruning. *Artificial intelligence*, 6(4):293–326.

Kocsis, L. and Szepesvári, C. (2006). Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer.

Kovařík, V., Schmid, M., Burch, N., Bowling, M., and Lisỳ, V. (2019). Rethinking formal models of partially observable multiagent decision making. *arXiv preprint arXiv:1906.11110*.

Kovarık, V., Seitz, D., Lisỳ, V., Rudolf, J., Sun, S., and Ha, K. (2020). Value functions for depth-limited solving in imperfect-information games. *arXiv preprint arXiv:1906.06412*.

Kovařík, V. and Lisý, V. (2019). Problems with the efg formalism: a solution attempt using observations. *arXiv preprint arXiv:1906.06291*.

Kuhn, H. (1953). Extensive games and the problem of information. kuhn hw, tucker aw, eds., contributions to the theory of games, vol ii, 193–216.

Kuhn, H. W. (1950). A simplified two-person poker. *Contributions to the Theory of Games*, 1:97–103.

Lanctot, M. (2013). *Monte Carlo Sampling and Regret Minimization for Equilibrium Computation and Decision-Making in Large Extensive Form Games*. PhD thesis, University of Alberta, University of Alberta, Computing Science, 116 St. and 85 Ave., Edmonton, Alberta T6G 2R3.

Lanctot, M., Lisy, V., and Bowling, M. (2014). Search in imperfect information games using online monte carlo counterfactual regret minimization. In *Workshops at the Twenty-Eighth AAAI Conference on Artificial Intelligence*.

Lanctot, M., Waugh, K., Zinkevich, M., and Bowling, M. H. (2009). Monte carlo sampling for regret minimization in extensive games. In *NIPS*, pages 1078–1086.

Lanctot, M., Zambaldi, V., Gruslys, A., Lazaridou, A., Tuyls, K., Pérolat, J., Silver, D., and Graepel, T. (2017). A unified game-theoretic approach to multiagent reinforcement learning. *arXiv preprint arXiv:1711.00832*.

LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *nature*, 521(7553):436–444.

Leslie, D. S. and Collins, E. J. (2006). Generalised weakened fictitious play. *Games and Economic Behavior*, 56(2):285–298.

Li, H., Hu, K., Ge, Z., Jiang, T., Qi, Y., and Song, L. (2018). Double neural counterfactual regret minimization. *arXiv preprint arXiv:1812.10607*.

Ling, C. K. and Brown, N. (2021). Safe search for stackelberg equilibria in extensive-form games. *arXiv preprint arXiv:2102.01775*.

Lisy, V. and Bowling, M. (2017). Eqilibrium approximation quality of current no-limit poker bots. In *Workshops at the Thirty-First AAAI Conference on Artificial Intelligence*.

Lisý, V., Lanctot, M., and Bowling, M. (2015). Online monte carlo counterfactual regret minimization for search in imperfect information games. In *Proceedings of the 2015 international conference on autonomous agents and multiagent systems*, pages 27–36.

Littlestone, N. and Warmuth, M. K. (1994). The weighted majority algorithm. *Information and computation*, 108(2):212–261.

Littman, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pages 157–163. Elsevier.

McAleer, S., Lanier, J., Baldi, P., and Fox, R. (2021). Xdo: A double oracle algorithm for extensive-form games. *arXiv preprint arXiv:2103.06426*.

McAleer, S., Lanier, J., Fox, R., and Baldi, P. (2020). Pipeline psro: A scalable approach for finding approximate nash equilibria in large games. *arXiv preprint arXiv:2006.08555*.

McCarthy, J. (1990). Chess as the drosophila of ai. In *Computers, chess, and cognition*, pages 227–237. Springer.

McMahan, H. B., Gordon, G. J., and Blum, A. (2003). Planning in the presence of cost functions controlled by an adversary. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 536–543.

Moravčík, M., Schmid, M., Burch, N., Lisỳ, V., Morrill, D., Bard, N., Davis, T., Waugh, K., Johanson, M., and Bowling, M. (2017). Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513.

Moravcik, M., Schmid, M., Ha, K., Hladik, M., and Gaukrodger, S. J. (2016). Refining subgames in large imperfect information games. In *Thirtieth AAAI Conference on Artificial Intelligence*.

Morgenstern, O. and Von Neumann, J. (1953). *Theory of games and economic behavior*. Princeton university press.

Muller, P., Omidshafiei, S., Rowland, M., Tuyls, K., Perolat, J., Liu, S., Hennes, D., Marris, L., Lanctot, M., Hughes, E., et al. (2019). A generalized training approach for multiagent learning. *arXiv preprint arXiv:1909.12823*.

Nash, J. (1951). Non-cooperative games. *Annals of mathematics*, pages 286–295.

Nash, J. F. et al. (1950). Equilibrium points in n-person games. *Proceedings of the national academy of sciences*, 36(1):48–49.

Nau, D. S. (1982). An investigation of the causes of pathology in games. *Artificial Intelligence*, 19(3):257–278.

Neumann, J. v. (1928). Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320.

Newell, A. and Simon, H. A. (1964). An example of human chess play in the light of chess playing programs. Technical report, CARNEGIE INST OF TECH PITTSBURGH PA.

Nisan, N., Roughgarden, T., Tardos, E., and Vazirani, V. V. (2007). *Algorithmic game theory.* Cambridge university press.

Osborne, M. J. and Rubinstein, A. (1994). *A course in game theory.* MIT press.

Owen, A. B. (2016). Monte carlo theory, methods and examples. 2013. *URL http://statweb. stanford. edu/~ owen/mc.*

Papadimitriou, C. H. (1994). On the complexity of the parity argument and other inefficient proofs of existence. *Journal of Computer and system Sciences*, 48(3):498–532.

Pascutto, G.-C. (2019). Leela zero. `https://zero.sjeng.org`. [Online; accessed 05-Dec-2020].

Perolat, J., Munos, R., Lespiau, J.-B., Omidshafiei, S., Rowland, M., Ortega, P., Burch, N., Anthony, T., Balduzzi, D., De Vylder, B., et al. (2020). From poincar\'e recurrence to convergence in imperfect information games: Finding equilibrium via regularization. *arXiv preprint arXiv:2002.08456.*

Piccione, M. and Rubinstein, A. (1997). On the interpretation of decision problems with imperfect recall. *Games and Economic Behavior*, 20(1):3–24.

Piccione, M., Rubinstein, A., et al. (1996). *The absent minded driver's paradox: Synthesis and responses.* Sackler Institute for Economic Studies.

Rapoport, A. and Chammah, A. M. (1966). The game of chicken. *American Behavioral Scientist*, 10(3):10–28.

Rapoport, A., Chammah, A. M., and Orwant, C. J. (1965). *Prisoner's dilemma: A study in conflict and cooperation*, volume 165. University of Michigan press.

Robinson, J. (1951). An iterative method of solving a game. *Annals of mathematics*, pages 296–301.

Samuel, A. L. (1959). Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, 3(3):210–229.

Sandholm, T. (2010). The state of solving large incomplete-information games, and application to poker. *Ai Magazine*, 31(4):13–32.

Schmid, M. (2013). Game theory and poker. Master's thesis, Charles University in Prague.

Schmid, M., Burch, N., Lanctot, M., Moravcik, M., Kadlec, R., and Bowling, M. (2019). Variance reduction in monte carlo counterfactual regret minimization (vr-mccfr) for extensive form games using baselines. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 2157–2164.

Schmid, M. and Moravcik, M. (2013). Equilibrium's action bound in extensive form games with many actions. In *Workshops at the Twenty-Seventh AAAI Conference on Artificial Intelligence.*

Schmid, M., Moravcik, M., and Hladik, M. (2014). Bounding the support size in extensive form games with imperfect information. In *Twenty-Eighth AAAI Conference on Artificial Intelligence.*

Schmid, M., Moravcik, M., Hladik, M., and Gaukrodger, S. J. (2015). Automatic public state space abstraction in imperfect information games. In *AAAI Workshop: Computer Poker and Imperfect Information.*

Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural networks*, 61:85–117.

Schnizlein, D., Bowling, M., and Szafron, D. (2009). Probabilistic state translation in extensive games with large action sets. In *Twenty-First International Joint Conference on Artificial Intelligence.*

Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., et al. (2020). Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609.

Seitz, D., Kovařík, V., Lisỳ, V., Rudolf, J., Sun, S., and Ha, K. (2019). Value functions for depth-limited solving in imperfect-information games beyond poker. *arXiv preprint arXiv:1906.06412.*

Serrino, J., Kleiman-Weiner, M., Parkes, D. C., and Tenenbaum, J. (2019). Finding friend and foe in multi-agent games. In *Advances in Neural Information Processing Systems*, pages 1251–1261.

Shannon, C. E. (1950). Xxii. programming a computer for playing chess. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 41(314):256–275.

Shapley, L. S. (1953). Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100.

Shoham, Y. and Leyton-Brown, K. (2008). *Multiagent systems: Algorithmic, game-theoretic, and logical foundations.* Cambridge University Press.

Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489.

Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al. (2017a). Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815.*

Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. (2017b). Mastering the game of go without human knowledge. *nature*, 550(7676):354–359.

Southey, F., Bowling, M., Larson, B., Piccione, C., Burch, N., Billings, D., and Rayner, C. (2005). Bayes' bluff: opponent modelling in poker. In *Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence*, pages 550–558.

Steinberger, E. (2019). Single deep counterfactual regret minimization. *arXiv preprint arXiv:1901.07621*.

Steinberger, E., Lerer, A., and Brown, N. (2020). Dream: Deep regret minimization with advantage baselines and model-free learning. *arXiv preprint arXiv:2006.10410*.

Sustr, M., Kovarik, V., and Lisy, V. (2018). Monte carlo continual resolving for online strategy computation in imperfect information games. *arXiv preprint arXiv:1812.07351*.

Šustr, M., Schmid, M., Moravčík, M., Burch, N., Lanctot, M., and Bowling, M. (2020). Sound search in imperfect information games. *arXiv preprint arXiv:2006.08740*.

Šustr, M., Schmid, M., Moravčík, M., Burch, N., Lanctot, M., and Bowling, M. (2021). Sound algorithms in imperfect information games. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1674–1676.

Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction.* MIT press.

Sutton, R. S., Barto, A. G., et al. (1998). *Introduction to reinforcement learning*, volume 135. MIT press Cambridge.

Syrgkanis, V., Agarwal, A., Luo, H., and Schapire, R. E. (2015). Fast convergence of regularized learning in games. *arXiv preprint arXiv:1507.00407*.

Tammelin, O. (2014). Solving large imperfect information games using cfr+. *arXiv preprint arXiv:1407.5042*.

Tammelin, O., Burch, N., Johanson, M., and Bowling, M. (2015). Solving heads-up limit texas hold'em. In *Twenty-fourth international joint conference on artificial intelligence*.

Tesauro, G. (1995). Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68.

Timbers, F., Lockhart, E., Schmid, M., Lanctot, M., and Bowling, M. (2020). Approximate exploitability: Learning a best response in large games. *arXiv preprint arXiv:2004.09677*.

Tucker, G., Bhupatiraju, S., Gu, S., Turner, R. E., Ghahramani, Z., and Levine, S. (2018). The mirage of action-dependent baselines in reinforcement learning. *arXiv preprint arXiv:1802.10031*.

Van der Genugten, B. (2000). A weakened form of fictitious play in two-person zero-sum games. *International Game Theory Review*, 2(04):307–328.

Vinyals, O., Babuschkin, I., Chung, J., Mathieu, M., Jaderberg, M., Czarnecki, W. M., Dudzik, A., Huang, A., Georgiev, P., Powell, R., et al. (2019a). Alphastar: Mastering the real-time strategy game starcraft ii. *DeepMind blog*, 2.

Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. (2019b). Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354.

Von Neumann, J. and Fréchet, M. (1953). Communication on the borel notes. *Econometrica: journal of the Econometric Society*, pages 124–127.

Von Neumann, J. and Morgenstern, O. (2007). *Theory of games and economic behavior (commemorative edition)*. Princeton university press.

Waugh, K., Morrill, D., Bagnell, J. A., and Bowling, M. (2015). Solving games with functional regret estimation. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.

Waugh, K., Schnizlein, D., Bowling, M. H., and Szafron, D. (2009). Abstraction pathologies in extensive games. In *AAMAS (2)*, pages 781–788.

White, M. and Bowling, M. H. (2009). Learning a value analysis tool for agent evaluation. In *IJCAI*, pages 1976–1981. Citeseer.

Whitehouse, D. (2014). *Monte Carlo tree search for games with hidden information and uncertainty*. PhD thesis, University of York.

Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256.

Wolpert, D. H. (1996). The lack of a priori distinctions between learning algorithms. *Neural computation*, 8(7):1341–1390.

Zarick, R., Pellegrino, B., Brown, N., and Banister, C. (2020). Unlocking the potential of deep counterfactual value networks. *arXiv preprint arXiv:2007.10442*.

Zhang, B. H. and Sandholm, T. (2020). Small nash equilibrium certificates in very large games. *arXiv preprint arXiv:2006.16387*.

Zinkevich, M. (2003). Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (icml-03)*, pages 928–936.

Zinkevich, M., Johanson, M., Bowling, M. H., and Piccione, C. (2007). Regret minimization in games with incomplete information. In *NIPS*, volume 7, page 1729.

# List of Figures

# List of Tables

# Glossary

**AIVAT** Variance reduction technique for evaluating agent's performance. 138

**alpha-beta pruning** Pruning trick that allows the minimax algorithm to potentially visit less nodes of the game tree. 32

**behavioral strategy** Strategy defined over all the information states in extensive form games. 57

**best responding sequence** Sequence of strategies where the agent's next strategy is best response against the the last strategy of the opponent. 34

**best response** Strategy maximizing the player's utility given the fixed policy of other players. 161

**common information set** Set of states closed under the consistency operation. 72

**consistency operation** Operation that given a state, returns consistent states. 72

**consistent states** Set of (opponent's) states consistent with a current state.. 72, 161

**constant-sum game** Strictly competitive game where the reward of both players always sum to a constant. 20, 162

**counterfactual best response** Best response having non-negative counterfactual regrets. 97

**Counterfactual regret minimization** Regret minimization in sequential decision making (i.e on a game tree). 94

**equilibrium selection problem** Selecting Nash equilibrium to play in multiplayer settings. 27, 70

**extensive form games** Formalism for imperfect information games. 56, 161

**factored observation stochastic games** Formalism for imperfect information games. 58, 60

**fictitious play** Variant of self-play where the agent's next strategy in the sequence is best response against the averaged strategy of the opponent. 88

**game tree** Formal model of sequential decision making in perfect information settings. 19

**game value** Unique value of the game when both players follow optimal policy. 24, 26

**Glasses** Simple variant of the graph chasing game. 54

# Acronyms

$\mathbb{NEQ}$ set of all Nash equilibrium strategy profiles. 24

**CFR** Counterfactual regret minimization. *Glossary:* Counterfactual regret minimization

# List of Publications

## Journal Papers

Moravcik Matej*, Martin Schmid*, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. "Deepstack: Expert-level artificial intelligence in heads-up no-limit poker." Science 356, no. 6337 (2017): 508-513.

Burch, N., Moravcik, M. and Schmid, M., 2019. Revisiting cfr+ and alternating updates. Journal of Artificial Intelligence Research, 64, pp.429-443.

## Conference Papers

Sustr, M., Schmid, M., Moravcik, M., Burch, N., Lanctot, M., & Bowling, M. (2020). "Sound search in imperfect information games" 20th International Conference on Autonomous Agents and Multiagent Systems

Davis, Trevor, Martin Schmid, and Michael Bowling. "Low-Variance and Zero-Variance Baselines for Extensive-Form Games." International Conference on Machine Learning. PMLR, 2020.

Schmid, Martin, et al. "Variance reduction in monte carlo counterfactual regret minimization (VR-MCCFR) for extensive form games using baselines." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. No. 01. 2019.

Burch, N., Schmid, M., Moravcik, M., Morill, D. and Bowling, M., 2018, April. Aivat: A new variance reduction technique for agent evaluation in imperfect information games. In Thirty-Second AAAI Conference on Artificial Intelligence.

Moravcik, M., Schmid, M., Ha, K., Hladik, M. and Gaukrodger, S.J., 2016, February. Refining subgames in large imperfect information games. In Thirtieth AAAI Conference on Artificial Intelligence.

Schmid, M., Moravcik, M. and Hladik, M., 2014, June. Bounding the support size in extensive form games with imperfect information. In Twenty-Eighth AAAI Conference on Artificial Intelligence.

## Workshop Papers

Schmid, Martin, et al. "Automatic public state space abstraction in imperfect information games." AAAI Workshop: Computer Poker and Imperfect Information. 2015.

---

[1]*Equal contribution, alphabetical order.

# ArXiv

Gruslys, Audrunas, et al. "The advantage regret-matching actor-critic." arXiv preprint arXiv:2008.12234 (2020).

Sokota, Samuel, et al. "Solving Common-Payoff Games with Approximate Policy Iteration." arXiv preprint arXiv:2101.04237 (2020).

Timbers, F., Lockhart, E., Schmid, M., Lanctot, M., Bowling, M. (2020). Approximate exploitability: Learning a best response in large games. arXiv preprint arXiv:2004.09677.

Kovarik, V., Schmid, M., Burch, N., Bowling, M. and Lisy, V., 2019. Rethinking formal models of partially observable multiagent decision making. arXiv preprint arXiv:1906.11110.