

Model-Free Neural Counterfactual Regret Minimization with Bootstrap Learning

Weiming Liu, Bin Li,* *Member, IEEE*, and Julian Togelius, *Member, IEEE*

Abstract—Counterfactual Regret Minimization (CFR) has achieved many fascinating results in solving large-scale Imperfect Information Games (IIGs). Neural network approximation CFR (neural CFR) is one of the promising techniques that can reduce computation and memory consumption by generalizing decision information between similar states. Current neural CFR algorithms have to approximate cumulative regrets. However, efficient and accurate approximation in a large-scale IIG is still a tough challenge. In this paper, a new CFR variant, Recursive CFR (ReCFR), is proposed. In ReCFR, Recursive Substitute Values (RSVs) are learned and used to replace cumulative regrets. It is proven that ReCFR can converge to a Nash equilibrium at a rate of $O(1/\sqrt{T})$. Based on ReCFR, a new model-free neural CFR with bootstrap learning, Neural ReCFR-B, is proposed. Due to the recursive and non-cumulative nature of RSVs, Neural ReCFR-B has lower-variance training targets than other neural CFRs. Experimental results show that Neural ReCFR-B is competitive with the state-of-the-art neural CFR algorithms at a much lower training cost.

Index Terms—Game Theory, Imperfect Information Games, Counterfactual Regret Minimization, Neural Networks.

I. INTRODUCTION

IMPERFECT Information Game (IIG) is a kind of challenging game, in which players can only obtain partial information. Conventional methods for Perfect Information Games (PIGs) [1], [2], [3], [4] generally do not play IIGs well [5], [6], [7]. Usually, IIGs are solved using equilibrium-finding algorithms [8], [9], [10], [11]. This paper focuses on two-player, zero-sum IIGs in extensive form, for example, heads-up no-limit Texas hold'em poker (HUNL). In recent years, Counterfactual Regret Minimization (CFR) algorithm [8] has been applied in human-level poker agents [12], [13], [14]. CFR is an iterative algorithm. It has been proven that the average strategy will converge to a Nash equilibrium if the zero-thresholded regret, $\max\{\max_a R_p^T(I, a), 0\}$, at every state increases sub-linearly [8]. Here $R_p^T(I, a) = \sum_{t=1}^T r^t(I, a)$ is the *cumulative regret*, and $r^t(I, a)$ is the *instantaneous regret*, will be defined later. To minimize the cumulative regret at every state, a CFR algorithm needs to choose a strategy at every iteration according to the regrets. Conventional tabular CFR visits all states at every iteration and stores the cumulative

regrets in a table. So it is computation expensive and memory intensive in large-scale IIGs.

To deal with large-scale IIGs, the technique “abstraction” [15] was proposed to cluster states and actions. However, abstraction techniques are domain-specific and depend on expert knowledge. Function approximation CFR (FCFR) [16], [17] is a kind of algorithm that approximates cumulative regrets with function approximators. Recently, some FCFRs with neural network approximators (neural CFR) have been proposed [18], [19], [20]. Thanks to the potential generalization ability of function approximators, in an FCFR, it can be considered that the states are clustered in an implicit space. In other words, FCFR may abstract a game automatically, and, hopefully, reducing the computation and memory consumption.

Since it is intractable to traverse the full game tree and approximate the cumulative regrets directly in a large-scale IIG, FCFR usually trains the approximator on surrogate training targets. As proven in [17], [18], the quality of the approximation directly impacts the quality of the output strategy. However, the training cost for approximating cumulative regrets could be high. Take a linear approximator with p parameters as an example. Assume the variance in training targets is σ^2 , and the data set size is m . It is known that the variance in approximated values on the fixed data set is $\frac{p}{m}\sigma^2$ [21]. Unfortunately, the variance in training targets in a zero-sum IIG is likely to be high due to the adversarial nature. So, as a compromise, the data set should be large. Also, the approximator should be trained for many epochs to prevent underfitting. Therefore, the cost for approximating cumulative regrets could be high. We will discuss more in Section III-D.

Furthermore, existing neural CFR algorithms usually require an exact simulator for model-based sampling (sampling multiple actions in a single state), in order to reduce the variance in training targets. To extend CFR to more IIGs with unknown models, it is necessary to develop model-free algorithms.

To reduce the training cost and use model-free sampling, it is desired to avoid approximating cumulative regrets. In this paper, a new CFR algorithm, Recursive CFR (ReCFR), is proposed, in which cumulative regrets are replaced by non-cumulative Recursive Substitute Values (RSVs) proposed in Warm CFR [22]. We prove that ReCFR converges to a Nash equilibrium at a rate of $O(1/\sqrt{T})$. We also prove that vanilla CFR and Full-width extensive-form fictitious play (XFP) [9] are special cases of ReCFR.

Based on ReCFR, a new neural CFR algorithm, Neural ReCFR with Bootstrapping (Neural ReCFR-B), is proposed. In the algorithm, two RSV networks (one for each player) are trained to approximate RSVs. The training targets for RSVs

Weiming Liu is with the School of Data Science, University of Science and Technology of China, Anhui, China. (e-mail: weiming@mail.ustc.edu.cn).

Bin Li (corresponding author) is with the School of Information Science and Technology, University of Science and Technology of China, Anhui, China. (e-mail: binli@ustc.edu.cn).

Julian Togelius is with the Department of Computer Science and Engineering, New York University, New York, NY 11201, USA (e-mail: julian@togelius.com).

The work is partially supported by the National Natural Science Foundation of China under grand No.U19B2044 and No.61836011.

can be estimated by bootstrap learning, which is an important benefit brought by the recursive nature of RSVs (cumulative regrets do not have this property). Since RSVs are non-cumulative, the training targets should have lower variance than that for cumulative regrets. Also, bootstrap learning can reduce the variance, at the cost of being a bit biased [23]. So, due to the lower variance training targets, Neural ReCFR-B may use model-free sampling and train the neural networks with fewer samples at every iteration than other neural CFR algorithms. Experimental results show that Neural ReCFR-B is competitive with the state-of-the-art neural CFR algorithms [18], [19], [20] at a much lower training cost.

The contributions of the paper are mainly in three aspects:

- RSVs and Warm CFR [22] are revisited, and a new exploitability bound is proven. Two new and intuitive lemmas (Lemma 2 and Lemma 3) are also proposed for analyzing CFR algorithms.
- A new CFR algorithm based on RSVs is presented, and it is proven to converge to a Nash equilibrium at a rate of $O(1/\sqrt{T})$. It is also shown that ReCFR can generalize both Fictitious Play [9] and vanilla CFR.
- A new model-free neural CFR algorithm based on ReCFR is presented, and it shows higher training efficiency than other algorithms on medium-size and large-size IIGs.

In the rest of the paper, related work is first discussed in Section II. Then, notations and background are given in Section III. In Section IV, RSVs and Warm CFR are revisited, and two new lemmas are given. In Section V and VI, ReCFR and Neural ReCFR-B are described in detail. Then, the experimental setup and the results of ReCFR and Neural ReCFR-B are shown in Section VII. Finally, the conclusion is drawn in Section VIII.

II. RELATED WORK

In recent years, many techniques have been proposed to improve the performance of CFR in large-scale IIGs [24], [6], [25], [12]. Based on these techniques, the state-of-the-art algorithms for HUNL are Libratus [13] and Pluribus [14].

Regression CFR [16] is the first FCFR. Deep CFR [18] and Double Neural CFR (DNCFR) [19] are two more recent neural CFRs. Deep CFR uses two reservoir memory buffers to store all past instantaneous regrets and uses two regret networks (one for each player) to approximate the mean values of the past regrets. However, due to the high variance training targets, the variance in approximated cumulative regrets could be high. Another neural CFR, DNCFR, sums the current-iterate instantaneous regrets and the previous approximation as the training targets. However, the approximation error in each iteration will accumulate. So, as the number of iterations goes to infinity, the bias is unbounded. The problem may be alleviated using Regret Matching+ algorithm [26] and model-based robust sampling [19], but there is no theoretical guarantee. Single Deep CFR [27] is a variant of Deep CFR, which does not need to approximate the average strategy. Instead, it keeps all regret networks of past iterations and computes an average strategy at the end. However, cumulative regrets are still approximated.

Neural Fictitious Self-Play (NFSP) [28] is a model-free algorithm based on full-width extensive-form fictitious play (XFP) [9]. NFSP uses DQN [3] to learn the best-response strategy for each player at each iteration and uses a strategy network to approximate the average best-response strategy. Since CFR has a better convergence guarantee than XFP, people may also be interested in model-free neural CFR algorithms. ARMAC [29] and DREAM [20] are both model-free neural CFR algorithms like Neural ReCFR-B. However, both of them approximate cumulative regrets. Specifically, ARMAC approximates cumulative regrets by replaying past strategies and regenerating past regrets, while DREAM uses a reservoir buffer to store all past regrets. Besides, they also train an additional state-action network to reduce the variance in training targets [30]. A critical difference of Neural ReCFR-B from these two methods is that it does not store or regenerate past regrets but learns RSVs only based on the average strategy.

Some other works have combined equilibrium-finding algorithms with gradient descending [31], [32]. They have achieved performance comparable to NFSP, and they seem to have the potential for solving large-scale IIGs. However, they have only been tested in small-scale games.

III. NOTATIONS AND BACKGROUND

For a two-player zero-sum IIG, the set of players is denoted by $P = \{1, 2\}$. The chance player, denoted by c , is introduced to take actions for random events. A history h is represented as a sequence of actions taken by all the players and the chance player. The root history is represented as an empty sequence, written as \emptyset . Denote the set of histories by H . For any non-terminal history $h \in H$, the set of legal actions at h is denoted by $A(h)$. The set of all actions is denoted by \mathcal{A} . The acting player at history h is denoted by $P(h)$, where $P(h) \in P \cup \{c\}$. In this paper, we assume that $P(\emptyset) \in P$ and will present the main results under the viewpoint of $p = P(\emptyset)$.¹ After player $P(h)$ takes an action $a \in A(h)$, the resulted history is denoted by $h \cdot a$. If there exists a sequence of actions from h to h' , then h' is a descendent of h , denoted by $h \sqsubset h'$. Let $h \sqsubseteq h'$ represent that $h \sqsubset h'$ or $h = h'$. We only consider depth-limited games in this paper. The set of terminal histories is denoted by Z . For any $z \in Z$, the *payoff* for player p is $u_p(z) \in \mathbb{R}$. As the game is two-player and zero-sum, we have $u_1(z) = -u_2(z)$. Define $Succ_p(h)$ as the set of the earliest reachable histories such that for any $h' \in Succ_p(h)$, $P(h') = p$ or $h' \in Z$. ‘‘Earliest reachable’’ means that $h \sqsubseteq h'$ and there is no h'' such that $h'' \sqsubset h'$.

In an imperfect information game, the histories of each player $p \in P$ are divided into *information sets* (infosets). The set of infosets of player p is denoted by \mathcal{I}_p . Let $\mathcal{I} = \mathcal{I}_1 \cup \mathcal{I}_2$. For any infoset $I \in \mathcal{I}_p$, all histories $h, h' \in I$ are indistinguishable to p , so $A(h) = A(h')$. Define $P(I) = P(h)$ and $A(I) = A(h)$ for any $h \in I$. Let $I(h)$ denote the infoset of history h , i.e., $I(h) = I$ for any

¹For the other player $p' \neq P(\emptyset)$ we can construct an equivalent surrogate game with $P(\emptyset) = p'$ by adding a player p' decision point at the beginning of the game with a single action [22].

$h \in I$. Besides, define $Succ_p(I, a)$ as the set of the earliest reachable infosets after action a has been taken. Formally, $Succ_p(I, a) = \{I(h') | h' \notin Z, h' \in Succ_p(h \cdot a), h \in I\}$. Let $Succ_p(I) = \bigcup_{a \in A(I)} Succ_p(I, a)$. The range of payoffs reachable from I is denoted by $\Delta(I)$. Formally, $\Delta(I) = \max_{z \in Z, h \in I: h \sqsubset z} u_P(I)(z) - \min_{z \in Z, h \in I: h \sqsubset z} u_P(I)(z)$. Let $\Delta = \max_{I \in \mathcal{I}} \Delta(I)$.

A strategy σ_p of a player $p \in P$ is a function that maps any infoset $I \in \mathcal{I}_p$ to a probability vector over $A(I)$. The set of strategies of player p is denoted by Σ_p . Given a strategy $\sigma_p \in \Sigma_p$, $\sigma_p(I)$ is the probability vector, and $\sigma_p(I, a)$ is the probability of choosing action a at infoset I . Since the histories in an infoset are indistinguishable, the strategy in each of them must be identical. So $\sigma_p(h) = \sigma_p(I)$ and $\sigma_p(h, a) = \sigma_p(I, a)$ for any $h \in I$ and $a \in A(I)$. The strategy of the other player (the opponent) is denoted by σ_{-p} . Denote a strategy profile by $\sigma = \langle \sigma_p, \sigma_{-p} \rangle$. The chance player's strategy, which is fixed and known to all the players, is denoted by $\sigma_c(h, a)$.

$\pi^\sigma(h) = \prod_{h' a \sqsubset h} \sigma_P(h')(h', a)$ is called the *reach* of h , which is the probability of reaching h when all the players act according to σ . $\pi_p^\sigma(h)$ is the contribution of p to this probability. $\pi_{-p}^\sigma(h)$ is the contribution of the opponent and the chance player. The probability of reaching h' from h is denoted by $\pi^\sigma(h, h')$. In this paper, we only consider *perfect recall* games. Therefore, we define $\pi_p^\sigma(I) = \pi_p^\sigma(h)$ for any $h \in I$. Accordingly, define the reach of the opponent and the chance player as $\pi_{-p}^\sigma(I) = \sum_{h \in I} \pi_{-p}^\sigma(h)$. The *expected payoff* of the game for player p is denoted by $u_p(\sigma_p, \sigma_{-p})$. Formally, $u_p(\sigma_p, \sigma_{-p}) = \sum_{z \in Z} \pi^\sigma(z) u_p(z)$. A *best response* $BR(\sigma_{-p})$ is a strategy of player p such that $BR(\sigma_{-p}) = \operatorname{argmax}_{\sigma'_p \in \Sigma_p} u_p(\sigma'_p, \sigma_{-p})$. A *Nash equilibrium* $\sigma^* = \langle \sigma_p^*, \sigma_{-p}^* \rangle$ is a strategy profile where every player plays a best response. The *exploitability* of a strategy σ_p is the distance to a Nash equilibrium, defined as $e(\sigma_p) = u_p(\sigma_p^*, BR(\sigma_p^*)) - u_p(\sigma_p, BR(\sigma_p))$. Define the total exploitability as $e(\sigma) = \sum_{p \in P} e(\sigma_p)$.

A. Counterfactual Regret Minimization (CFR)

CFR is an iterative algorithm for two-player zero-sum IIGs. It computes a strategy profile at every iteration using Regret Matching (RM) algorithm [8] according to the cumulative regrets of infosets. Then, the full game tree is traversed, and the cumulative regrets and average strategy are updated according to the strategy. It has been proven that the total exploitability of the average strategy is bounded by $O(1/\sqrt{T})$ after T iterations of CFR [8] are played.

Let σ^t be the strategy at iteration t . The *counterfactual value* of action a at infoset I is defined as

$$v_p^{\sigma^t}(I, a) = \sum_{h \in I} \sum_{z \in Z: h \sqsubset z} \pi_{-p}^{\sigma^t}(h) \pi^{\sigma^t}(h \cdot a, z) u_p(z). \quad (1)$$

$v_p^{\sigma^t}(I) = \sum_{a \in A(I)} \sigma_p^t(I, a) v_p^{\sigma^t}(I, a)$ is the counterfactual value of infoset I . Counterfactual values can also be defined recursively: at each infoset I ,

$$v_p^{\sigma^t}(I, a) = \sum_{h \in I} \sum_{z \in Z: z \in Succ_p(h \cdot a)} \pi_{-p}^{\sigma^t}(z) u_p(z) + \sum_{I' \in Succ_p(I, a)} v_p^{\sigma^t}(I'). \quad (2)$$

This equation has been used in some existing literature [8], [22]. A proof is also provided in Appendix A.² Let $v_p^{\sigma^t} = v_p^{\sigma^t}(I(\emptyset))$, then, $u_p(\sigma_p^t, \sigma_{-p}^t) = v_p^{\sigma^t}$. In the rest of the paper, we mainly use $v_p^{\sigma^t}$ and $v_p^{\langle \sigma_p^t, \sigma_{-p}^t \rangle}$ to denote the expected payoff. The *instantaneous regret* is defined as $r_p^t(I, a) = v_p^{\sigma^t}(I, a) - v_p^{\sigma^t}(I)$. The *cumulative counterfactual regret (cumulative regret)* of action a at I is

$$R_p^T(I, a) = \sum_{t=1}^T v_p^{\sigma^t}(I, a) - \sum_{t=1}^T v_p^{\sigma^t}(I). \quad (3)$$

The *average strategy* $\bar{\sigma}_p^T(I, a)$ is computed³ according to

$$\bar{\sigma}_p^T(I, a) = \frac{\sum_{t=1}^T \pi_p^{\sigma^t}(I) \sigma_p^t(I, a)}{\sum_{t=1}^T \pi_p^{\sigma^t}(I)}. \quad (4)$$

The *total regret* of player p after T iterations is

$$R_p^T = \max_{\sigma'_p \in \Sigma_p} \sum_{t=1}^T v_p^{\langle \sigma'_p, \sigma_{-p}^t \rangle} - \sum_{t=1}^T v_p^{\sigma^t}. \quad (5)$$

It has been proven in [8] that $R_p^T \leq \sum_{I \in \mathcal{I}_p} \max_a (R_p^T(I, a))_+$, where $(\cdot)_+ = \max\{\cdot, 0\}$. Therefore, the total regret R_p^T can be minimized by minimizing $\max_a (R_p^T(I, a))_+$ at every infoset, using, e.g., RM. At each infoset, RM computes the next iteration strategy according to

$$\sigma_p^{t+1}(I, a) = \frac{(R_p^t(I, a))_+}{\sum_{a' \in A(I)} (R_p^t(I, a'))_+}. \quad (6)$$

If $\sum_{a'} (R_p^t(I, a'))_+ = 0$, the action with the highest cumulative regret is assigned with probability 1 [18]. According to [22], we have the following lemma.

Lemma 1. [22] *After T iterations of CFR are played, for any infoset I , $\sum_{a \in A(I)} (R_p^T(I, a))_+^2 \leq \pi_{-p}^{\bar{\sigma}^T}(I) \Delta^2(I) |A(I)| T$.*

Finally, it is well known that, in a two-player zero-sum game, the total exploitability equals the average total regret, i.e., $e(\bar{\sigma}^T) = \frac{1}{T} \sum_{p \in P} R_p^T$. This can be seen by the definitions. Therefore, $e(\bar{\sigma}^T) = O(1/\sqrt{T})$ in CFR [8].

B. Recursive Substitute Values (RSVs) and Warm CFR

RSVs were first proposed in Warm CFR [22]. Given an arbitrary strategy σ , Warm CFR initializes the cumulative regrets with substitute regrets, hoping the CFR algorithm can converge faster. The initialization is done by treating the strategy σ as the average strategy after T iterations of CFR are run. To some extent, the initialization is equivalent to recovering the cumulative regrets. Let $v_p^{I\sigma}(I, a)$ and $v_p^{I\sigma}(I)$ be the RSVs that recover the true average counterfactual values: $\frac{1}{T} \sum_{t=1}^T v_p^{\sigma^t}(I, a)$ and $\frac{1}{T} \sum_{t=1}^T v_p^{\sigma^t}(I)$. Then, the *substitute cumulative counterfactual regret (substitute regret)*,

$$R_p^{IT}(I, a) = T (v_p^{I\sigma}(I, a) - v_p^{I\sigma}(I)), \quad (7)$$

should recover $R_p^T(I, a)$ at every infoset. According to the definition of counterfactual values and Lemma 1, the RSVs

²Appendix: <https://arxiv.org/abs/2012.01870>

³We assume $\bar{\sigma}^t(I, a) > 0$ (and thus $\pi_{-p}^{\bar{\sigma}^t}(I) > 0$) for any $I \in \mathcal{I}$ and $t \geq 1$. This is true if, e.g., $\sigma^1(I, a) = \frac{1}{|A(I)|} > 0$ for every $I \in \mathcal{I}$.

should satisfy the two constraints below if they are the true average counterfactual values:

$$v_p^{\prime\sigma}(I, a) = \sum_{h \in I} \sum_{z \in Z: z \in \text{Succ}_p(h, a)} \pi_{-p}^\sigma(z) u_p(z) + \sum_{I' \in \text{Succ}_p(I, a)} v_p^{\prime\sigma}(I'), \quad (8)$$

$$\sum_{a \in A(I)} (v_p^{\prime\sigma}(I, a) - v_p^{\prime\sigma}(I))_+^2 \leq \frac{\pi_{-p}^\sigma(I) \Delta^2(I) |A(I)|}{T}. \quad (9)$$

There could be many RSV profiles that fulfill the two constraints above. However, we can choose the RSVs for each info set recursively: at info set I , $v_p^{\prime\sigma}(I, a)$ is computed according to the immediate payoffs and the RSVs of its earliest reachable info sets, while $v_p^{\prime\sigma}(I)$ is chosen to fulfill (9).⁴ Define the *substitute expected payoff* as

$$v_p^{\prime\sigma} = v_p^{\prime\sigma}(I(\emptyset)). \quad (10)$$

In Warm CFR, another T' CFR iterations are run based on the substitute regrets. Specifically, define

$$R_p^{\prime T, T'}(I, a) = R_p^{\prime T}(I, a) + \sum_{t'=1}^{T'} r_p^{\sigma^{t'}}(I, a). \quad (11)$$

The strategy $\sigma_p^{t'}$ at iteration t' is computed using RM according to $R_p^{\prime T, t'-1}(I, a)$. Define the average strategy as $\bar{\sigma}_p^{\prime T, T'} = \frac{1}{T+T'}(T\sigma_p + T'\bar{\sigma}_p^{\prime T'})$, where $\bar{\sigma}_p^{\prime T'}$ is the average strategy of the T' iterations. It is proven in [22] that

$$\sum_{a \in A(I)} \left(R_p^{\prime T, T'}(I, a) \right)_+^2 \leq \pi_{-p}^{\bar{\sigma}_p^{\prime T, T'}}(I) \Delta^2(I) |A(I)| (T + T'). \quad (12)$$

As a result, the average strategy will converge to a Nash equilibrium at a rate of $O(1/\sqrt{T+T'})$ under specific condition (Assumption 1, will be discussed in the next section).

C. Monte Carlo CFR

Monte Carlo CFR (MCCFR) [33] is a variant of CFR that only traverses parts of the game tree at every iteration and estimates the instantaneous regrets by a Monte Carlo (MC) method. At iteration t , suppose a subset Q of the full game tree is sampled, the sampled counterfactual value at (I, a) is defined as

$$\tilde{v}_p^{\sigma^t}(I, a) = \sum_{z \in Q \cap Z} \frac{1}{q(z)} \pi_{-p}^{\sigma^t}(h) \pi^{\sigma^t}(h, a, z) u_p(z), \quad (13)$$

where $q(z)$ is the probability of sampling z . Define $\tilde{v}_p^{\sigma^t}(I) = \sum_{a \in A(I)} \sigma_p^{\sigma^t}(I, a) \tilde{v}_p^{\sigma^t}(I, a)$. The sampled instantaneous regret is $\tilde{r}_p^{\sigma^t}(I, a) = \tilde{v}_p^{\sigma^t}(I, a) - \tilde{v}_p^{\sigma^t}(I)$. The sampled cumulative regret is $\tilde{R}_p^T(I, a) = \sum_{t=1}^T \tilde{r}_p^{\sigma^t}(I, a)$, which is an unbiased estimator of the cumulative regret. There are many MCCFR algorithms [33], [34], [35], of which External Sampling CFR (ESCFR) and Outcome Sampling CFR (OSCFR) proposed in [33] are the most common.

⁴Because of the non-smooth zero-thresholding operator, $v_p^{\prime\sigma}(I)$ can not be presented in an explicit form. More details are provided in Appendix A.

D. Function Approximation of Cumulative Regrets

For a two-player IIG, FCFR is usually combined with MCCFR and approximates the MC sampled cumulative regrets by training on surrogate training targets [18], [19], [20]. Denote the training target at (I, a) by $\hat{R}^T(I, a)$ and assume $\hat{R}^T(I, a) = \tilde{R}^T(I, a) + \epsilon$, where $\tilde{R}^T(I, a)$ is the MC sampled cumulative regret and ϵ is a random variable at the info set. The variance in training targets is defined as $\text{var}(\hat{R}^T) = \mathbb{E}_{I, a, \epsilon} (\hat{R}^T(I, a) - \mathbb{E}_\epsilon \hat{R}^T(I, a))^2$. The expected error is $Err = \mathbb{E}_{I, a, D} (f^T(I, a|D) - \tilde{R}^T(I, a))^2$. According to the famous bias-variance decomposition, $Err = (\text{bias}(f^T))^2 + \text{var}(f^T)$, where $(\text{bias}(f^T))^2$ and $\text{var}(f^T)$ are the bias and variance in approximated regrets, respectively. Specifically, $(\text{bias}(f^T))^2 = \mathbb{E}_{I, a} (\mathbb{E}_D f^T(I, a|D) - \tilde{R}^T(I, a))^2$ and $\text{var}(f^T) = \mathbb{E}_{I, a, D} (f^T(I, a|D) - \mathbb{E}_D f^T(I, a|D))^2$. The bias could be low if the training targets are unbiased estimators. However, as we mentioned before, the variance in approximated cumulative regrets is affected by the variance in training targets and data set size.

As discussed in existing literature [17], [18], [19], there are at least two methods with different training targets for approximating sampled cumulative regrets. Note that a regret $\tilde{r}^{\sigma^t}(I, a), t = 1, \dots, T$ sampled according to MCCFR is an unbiased estimator of a kind of average regret at info set I . So, the first method, which has been used in Deep CFR and DREAM, is to train on a data set $D = \{(I, \tilde{r}^{\sigma^t}(I, \cdot)) | x^t(I) = 1, t = 1, \dots, T\}$, where $x^t(I)$ is an indicator variable that is 1 if and only if I is sampled at iteration t , and it is equal to 0 otherwise. Due to the adversarial nature of zero-sum IIGs, the sampled instantaneous regrets from different iterations are non-IID, and their variance tends to be high. To reduce the variance in training targets, model-based sampling is used in Deep CFR, and a variance reduction technique [30] is used in DREAM. A large regret memory buffer is also used. Besides, considering that the one-to-many mapping $(I, a) \mapsto \tilde{r}^{\sigma^t}(I, a)$ is complex, they train the approximators for many epochs to prevent underfitting.

The second method, which was first proposed in [17] and has been used in DNCFR [19], is to estimate $\tilde{R}^T(I, a)$ at iteration T with $f^{T-1}(I, a) + \tilde{r}^{\sigma^T}(I, a)$ and train the approximator on a data set $D = \{(I, f^{T-1}(I, \cdot) + \tilde{r}^{\sigma^T}(I, \cdot)) | I \in x^T(I)\}$. So, this is a method bootstrapping on the last approximation. The variance in training targets should be low since variance only comes from the current-iterate instantaneous regrets. However, the targets are *biased* estimators of the sampled cumulative regrets, and the approximation error can *accumulate* over iterations. In other words, the bias in approximated cumulative regrets is unbounded as the number of iterations goes to infinity. As a compromise, DNCFR trains the approximator for many epochs at every iteration. In conclusion, the cost of approximating cumulative regrets is high.

IV. REVISING RSVs AND WARM CFR

First, we would like to present a new and intuitive lemma for analyzing CFR algorithms.

Lemma 2. For any strategy profile $\sigma = \langle \sigma_p, \sigma_{-p} \rangle$, and another strategy σ'_p of player p , we have $v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - v_p^\sigma = \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'_p}(I) (v_p^\sigma(I, a) - v_p^\sigma(I)) \sigma'_p(I, a)$.

The lemma utilizes the recursive definition of counterfactual values (Equation (2)). Lemma 2 shows that the difference between expected payoffs can be completely represented by instantaneous regrets. As an application, the famous inequality $R_p^T \leq \sum_{I \in \mathcal{I}_p} \max_a (R_p^T(I, a))_+$ [8] is immediately recovered because $\sum_{t=1}^T v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - \sum_{t=1}^T v_p^\sigma = \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'_p}(I) \left(\sum_{t=1}^T r_p^t(I, a) \right) \sigma'_p(I, a)$.

Since RSVs have Equation (8) similar to (2), we have Lemma 3. Both lemmas are proven in Appendix A.

Lemma 3. For any strategy $\sigma = \langle \sigma_p, \sigma_{-p} \rangle$, another strategy σ'_p of player p , and arbitrary $v_p^{\sigma'}(I)$ at all infosets, compute $v_p^{\sigma'}(I, a)$ according to (8), then, $v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - v_p^\sigma = \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) (v_p^{\sigma'}(I, a) - v_p^\sigma(I)) \sigma'_p(I, a)$.

Note that Lemma 3 does not rely on how $v_p^{\sigma'}(I)$ is chosen at each infoset. Now, let us go back to the setting of Warm CFR. Assume the initial strategy σ is an average strategy $\bar{\sigma}^T$ generated by an arbitrary iterative algorithm, and another T' iterations of CFR are run based on the substitute regrets. As a result of Lemma 2 and 3, the total regret of the $T + T'$ iterations is

$$\begin{aligned} R_p^{T+T'} &= \max_{\sigma'_p \in \Sigma_p} \sum_{t=1}^{T+T'} v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - \sum_{t=1}^{T+T'} v_p^{\sigma^t} \\ &\leq \underbrace{\left(T v_p^\sigma - \sum_{t=1}^T v_p^{\sigma^t} \right)}_{\text{Term}_1} + \underbrace{\sum_{I \in \mathcal{I}_p} \max_a (R_p^{T, T'}(I, a))_+}_{\text{Term}_2}. \end{aligned} \quad (14)$$

A proof is provided in Appendix A. In the equation, the total regret of the $T + T'$ iterations is *decomposed* to two terms related to the RSVs. Specifically, **Term**₁ measures the difference between the substitute expected payoff and the true expected payoffs in the first T iterations, while **Term**₂ is the sum of the substitute regrets.

Assumption 1. [22] $v_1^\sigma + v_2^\sigma \leq 0$.

As shown in [22], if Assumption 1 is true, the total regret will be bounded by the sum of the substitute regrets, i.e., $\sum_{p \in P} R_p^{T+T'} \leq \sum_{p \in P} \sum_{I \in \mathcal{I}_p} \max_a (R_p^{T, T'}(I, a))_+$. Note that $\sum_{t=1}^T v_p^{\sigma^t}$ in **Term**₁ in (14) is canceled out as $v_1^\sigma + v_2^\sigma = 0$. Consequently, according to (12), we have $\epsilon(\bar{\sigma}^{T, T'}) = \frac{1}{T+T'} \sum_{p \in P} R_p^{T+T'} = O(1/\sqrt{T+T'})$. However, is Assumption 1 necessary for convergence? The answer may be no. Note that **Term**₁ in (14) does not depend on T' , and it may be bounded by the total regret of the initial T iterations. Based on these observations, we have the following theorem.

Theorem 1. Given an arbitrary strategy σ and $T > 0$, choose the RSVs according to (8) and (9) and assume $v_p^\sigma(I) \leq \max_a v_p^\sigma(I, a)$ at every infoset. If another T' iterations of CFR are run based on the substitute regrets, then, $\epsilon(\bar{\sigma}^{T, T'}) \leq \frac{T\epsilon(\sigma)}{T+T'} + \frac{1}{\sqrt{T+T'}} \sum_{I \in \mathcal{I}_p} \sqrt{\pi_{-p}^{\sigma^{T, T'}}(I) \Delta(I) \sqrt{|A(I)|}}$.

The proof is given in Appendix A. As we can see, when T' is sufficiently large (e.g., $T' \geq (T\epsilon(\sigma))^2 - T$), the first term $T\epsilon(\sigma)/(T+T') = O(1/\sqrt{T+T'})$, and thus $\epsilon(\bar{\sigma}^{T, T'}) = O(1/\sqrt{T+T'})$. For example, if the initial strategy is an average strategy of a CFR algorithm, we have $\epsilon(\sigma) = O(1/\sqrt{T})$ and $e(\sigma^{T, T'}) = O(\sqrt{T}/(T+T') + 1/\sqrt{T+T'}) = O(1/\sqrt{T+T'})$ for any $T' > 0$. Note that Theorem 1 does not rely on Assumption 1.

In this section, we relax Assumption 1 in Warm CFR to a trivial constraint: $v_p^{\sigma'}(I) \leq \max_a v_p^\sigma(I, a)$ for all infosets. Following the idea, we can prove the convergence of ReCFR given in the next section.

V. A NEW CFR ALGORITHM

In this section, a new CFR algorithm, named Recursive CFR (ReCFR), is presented. ReCFR is based on RSVs proposed in [22]. Instead of only computing the RSVs at the beginning for warm starting, ReCFR discards the cumulative regrets and replaces them with the substitute regrets at every iteration. Therefore, the cumulative regrets are never tracked. The pseudocode of ReCFR is given in Algorithm 1.

Algorithm 1 Recursive CFR

- 1: Input: maximum iterations T .
 - 2: **for** iteration $t = 1$ to T **do**
 - 3: **for** player $p \in P$ **do**
 - 4: **for** infoset $I \in \mathcal{I}_p$ **do**
 - 5: $\sigma_p^t(I, a) \leftarrow \begin{cases} \text{RM}(R_p^{t-1}(I, a)), & t > 1, \\ \frac{1}{|A(I)|}, & t = 1. \end{cases}$
 - 6: Update $\bar{\sigma}_p^t(I, a)$. ▷ (4)
 - 7: Compute $v_p^{\bar{\sigma}^t}(I, a)$ and $v_p^{\bar{\sigma}^t}(I)$. ▷ (8), (15)
 - 8: $R_p^t(I, a) \leftarrow t(v_p^{\bar{\sigma}^t}(I, a) - v_p^{\bar{\sigma}^t}(I))$.
-

Specifically, at iteration t , the RSV $v_p^{\bar{\sigma}^t}(I, a)$ is computed according to (8). However, we choose $v_p^{\bar{\sigma}^t}(I)$ at each infoset according to the constraint

$$\sum_{a \in A(I)} \left(t v_p^{\bar{\sigma}^t}(I, a) - t v_p^{\bar{\sigma}^t}(I) \right)_+^2 = \lambda_p^t(I), \quad (15)$$

where $\lambda_p^t(I) \in [0, \infty)$ is a parameter, will be set later. Note that when $\lambda_p^t(I) > 0$, the $v_p^{\bar{\sigma}^t}(I)$ fulfilling the constraint exists and is unique. Moreover, $v_p^{\bar{\sigma}^t}(I) < \max_a v_p^{\bar{\sigma}^t}(I, a)$. When $\lambda_p^t(I) = 0$, we force that $v_p^{\bar{\sigma}^t}(I) = \max_a v_p^{\bar{\sigma}^t}(I, a)$. The substitute regret $R_p^t(I, a)$ is

$$R_p^t(I, a) = t(v_p^{\bar{\sigma}^t}(I, a) - v_p^{\bar{\sigma}^t}(I)). \quad (16)$$

The strategy at iteration $t + 1$ is computed using RM:

$$\sigma_p^{t+1}(I, a) = \frac{(R_p^t(I, a))_+}{\sum_{a' \in A(I)} (R_p^t(I, a'))_+}. \quad (17)$$

When $\sum_{a' \in A(I)} (R_p^t(I, a'))_+ = 0$, the action with the maximal $R_p^t(I, a)$ is assigned with probability 1. We set $\sigma_p^1(I, a) = 1/|A(I)|$. The average strategy is computed according to (4). Since solving (15) requires a linear search, the cost of ReCFR

at every iteration is $O(|\mathcal{I}||\mathcal{A}|^2)$, which is worse than vanilla CFR ($O(|\mathcal{I}||\mathcal{A}|)$).

Note that when $\lambda_p^t(I) = 0$ at all infosets, we have $v_p^{(\sigma_p^{t+1}, \bar{\sigma}_{-p}^t)} = \max_{\sigma_p'} v_p^{(\sigma_p', \bar{\sigma}_{-p}^t)}$. In other words, σ_p^{t+1} is a best response to $\bar{\sigma}_{-p}^t$. According to the definition of XFP [9], we have proposition 1. Similarly, when $\lambda_p^t(I) = \sum_a (R^t(I, a))_+^2$ at each infoset, vanilla CFR is recovered. The proofs for the propositions are given in Appendix B.

Proposition 1. *ReCFR is equivalent to XFP if $\lambda_p^t(I) = 0$ at every infoset.*

Proposition 2. *ReCFR is equivalent to vanilla CFR if $\lambda_p^t(I) = \sum_a (R^t(I, a))_+^2$ at each infoset.*

A. Properties of Recursive CFR

According to (14), the total regret of ReCFR at iteration T can also be decomposed as

$$R_p^T \leq \underbrace{\left(T v_p^{i\bar{\sigma}^T} - \sum_{t=1}^T v_p^{\sigma^t} \right)}_{\text{Term}_1} + \underbrace{\sum_{I \in \mathcal{I}_p} \max_a (R_p^{iT}(I, a))_+}_{\text{Term}_2}. \quad (18)$$

Therefore, if Assumption 1 is true at iteration T , **Term**₁ can be canceled out when summing over R_p^T of the players, and thus $\epsilon(\bar{\sigma}^T) \leq \frac{1}{T} \sum_{p \in P} \sum_{I \in \mathcal{I}_p} \sqrt{\lambda_p^T(I)}$ (remember that $\epsilon(\bar{\sigma}^T) = \frac{1}{T} \sum_{p \in P} R_p^T$). However, it is *non-trivial* to ensure both Assumption 1 and $\lambda_p^T(I) = O(T)$ as $T \rightarrow \infty$.⁵ In other words, we may not be able to guarantee that $\epsilon(\bar{\sigma}^T) = O(1/\sqrt{T})$. Fortunately, similar to Theorem 1, Assumption 1 is unnecessary as long as both terms in (18) are bounded by $O(\sqrt{T})$. Specially, **Term**₁ equals

$$\sum_{t=1}^T \left(t v_p^{i\bar{\sigma}^t} - (t-1) v_p^{i\bar{\sigma}^{t-1}} - v_p^{\sigma^t} \right). \quad (19)$$

Thanks to Lemma 2 and Lemma 3, for $t \geq 1$, we have

$$\begin{aligned} & t v_p^{i\bar{\sigma}^t} - (t-1) v_p^{i\bar{\sigma}^{t-1}} - v_p^{\sigma^t} \\ &= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma^{t+1}}(I) g_p^{it}(I, a) \sigma_p^{t+1}(I, a), \end{aligned} \quad (20)$$

where $g_p^{it}(I, a) = \left(R_p^{t-1}(I, a) + r_p^{\sigma^t}(I, a) \right) - R_p^{it}(I, a)$. We set $R_p^0(I, a) = 0$ and $v_p^{i\bar{\sigma}^0} = 0$ as they are irrelevant to the algorithm. Here both $R_p^{t-1}(I, a) + r_p^{\sigma^t}(I, a)$ and $R_p^{it}(I, a)$ are the substitute regrets of infoset I at iteration t , according to (11). So it is potential that the difference is negligible. Based on the above analysis, Theorem 2 is obtained. The proof is given in Appendix B, where (20) is also proven.

Theorem 2. *After T iterations of ReCFR are played, for each player $p \in P$, if $\lambda_p^t(I) > 0$ at every infoset, then, $R_p^T \leq \sum_{t=1}^T \sum_{I \in \mathcal{I}_p} \frac{1}{2\sqrt{\lambda_p^t(I)}} \left(\lambda_p^{t-1}(I) - \lambda_p^t(I) + \sum_a (r_p^{\sigma^t}(I, a))^2 \right)_+ + \sum_{I \in \mathcal{I}_p} \sqrt{\lambda_p^T(I)}$.*

⁵Corollary 1 in [22] does not apply. It assumes T iterations of CFR were played (so Lemma 1 holds) when computing the RSVs, which is not true in ReCFR.

Theorem 2 shows that the two terms in (18) are bounded by the two terms in the theorem, respectively. When $\lambda_p^{t-1}(I) - \lambda_p^t(I) + \sum_a (r_p^{\sigma^t}(I, a))^2 \leq 0$, **Term**₁ in (18) will be bounded by zero, and thus $\epsilon(\bar{\sigma}^T) \leq \frac{1}{T} \sum_{p \in P} \sum_{I \in \mathcal{I}_p} \sqrt{\lambda_p^T(I)}$, as shown in Corollary 1.

Corollary 1. *If $\lambda_p^t(I) = \pi_{-p}^{\bar{\sigma}^t}(I) \Delta^2(I) |A(I)| t$ at each infoset, then, $\epsilon(\bar{\sigma}^T) \leq \frac{1}{\sqrt{T}} \sum_{p \in P} \sum_{I \in \mathcal{I}_p} \sqrt{\pi_{-p}^{\bar{\sigma}^T}(I) \Delta(I) \sqrt{|A(I)|}}$.*

This corollary is based on an inequality [22]:

$$\sum_a (r_p^{\sigma^t}(I, a))^2 \leq \pi_{-p}^{\sigma^t}(I) \Delta^2(I) |A(I)|. \quad (21)$$

Note that $\pi_{-p}^{\bar{\sigma}^t}(I) t = \sum_{k=1}^t \pi_{-p}^{\sigma^k}(I)$. As we can see, $\epsilon(\bar{\sigma}^T) = O(1/\sqrt{T})$. However, setting $\lambda_p^t(I)$ according to the corollary would perform poorly because the inequality in (21) could be loose. In other words, $\lambda_p^t(I)$ could be further reduced. Actually, we can choose $\lambda_p^t(I)$ in a broad range. According to Theorem 2, the total regret is bounded by $O(\sqrt{T})$ as long as $\lambda_p^{t-1}(I) - \lambda_p^t(I)$ is upper bounded and $\lambda_p^t(I) = \Theta(t)$.

Corollary 2. *If $\lambda_p^t(I) = \lambda \pi_{-p}^{\bar{\sigma}^t}(I) \Delta^2(I) |A(I)| t$, then $\epsilon(\bar{\sigma}^T) \leq \left(\frac{1}{\sqrt{\lambda}} + \sqrt{\lambda} \right) \frac{1}{\sqrt{T}} \sum_{p \in P} \sum_{I \in \mathcal{I}_p} \sqrt{\pi_{-p}^{\bar{\sigma}^T}(I) \Delta(I) \sqrt{|A(I)|}}$.*

The proofs for both corollaries are given in Appendix B. In Corollary 2, $\lambda \in (0, \infty)$ is a hyper-parameter. As we can see, the optimal exploitability bound is achieved when $\lambda = 1$. Although Corollary 2 provides a bound worse than that in Corollary 1, it allows λ to be chosen in $(0, \infty)$.

Note that Theorem 2 does not apply when $\lambda_p^t(I) = 0$ at some infosets. Therefore, the new theoretical results do not apply to CFR and XFP. Actually, when ReCFR is equivalent to CFR, **Term**₁ in (18) equals zero, and the known inequality $R_p^T \leq \sum_{I \in \mathcal{I}_p} \max_a (R_p^T(I, a))_+$ is recovered. When ReCFR is equivalent to XFP, **Term**₂ equals zero, and **Term**₁ is exactly the definition of R_p^T given in (5). So, Equation (18) does not provide new information for CFR or XFP. However, we can consider ReCFR a method that generalizes CFR and XFP.

B. Adapting the hyper-parameter

Corollary 2 suggests that the optimal λ should be 1. However, empirical results in Figure 1 show that the hyper-parameter could be much smaller. The reason may be that the bound for **Term**₁ in (18) is too loose. So, a smaller λ is needed to balance the two terms. Although Assumption 1 is not required in ReCFR, satisfying it may make the algorithm behave like a CFR. In this paper, we propose to use a simple adaptive algorithm to maintain that $v_1^{i\bar{\sigma}^t} + v_2^{i\bar{\sigma}^t} = 0$ loosely. Note that increasing $\lambda_p^t(I)$ at any infoset could reduce $v_p^{i\bar{\sigma}^t}$. So, we check $v_1^{i\bar{\sigma}^t} + v_2^{i\bar{\sigma}^t}$ at every iteration. If it is greater than 0, we increase λ : $\lambda = \beta_{amp} \lambda$ with $\beta_{amp} > 1$. Otherwise, λ is reduced: $\lambda = \beta_{damp} \lambda$ with $\beta_{damp} < 1$.

VI. A NEW NEURAL NETWORK APPROXIMATION CFR

This section describes Neural ReCFR with bootstrapping (Neural ReCFR-B) in detail. At each iteration, Neural ReCFR-B approximates the RSVs instead of the cumulative regrets.

The RSVs are computed according to (8) and (15). According to (8), $v_p^{\bar{\sigma}^t}(I, a)$ is scaled by the reach of the opponent and the chance player. Dividing it by $\pi_{-p}^{\bar{\sigma}^t}(I)$ will unify the ranges of the RSVs, which is helpful for neural approximation. Let $u_p^{\bar{\sigma}^t}(I, a) = v_p^{\bar{\sigma}^t}(I, a)/\pi_{-p}^{\bar{\sigma}^t}(I)$ and $u_p^{\sigma^t}(I) = v_p^{\sigma^t}(I)/\pi_{-p}^{\sigma^t}(I)$. In this paper, we choose to approximate $u_p^{\bar{\sigma}^t}(I, a)$ instead of $v_p^{\bar{\sigma}^t}(I, a)$. Put $u_p^{\bar{\sigma}^t}(I, a)$ into (15), then, $u_p^{\sigma^t}(I)$ is required to fulfill

$$\sum_{a \in A(I)} \left(u_p^{\bar{\sigma}^t}(I, a) - u_p^{\sigma^t}(I) \right)_+^2 = \beta_p^t(I), \quad (22)$$

where $\beta_p^t(I) = \lambda_p^t(I)/(\pi_{-p}^{\bar{\sigma}^t}(I)t)^2$. According to Corollary 2, we set $\beta_p^t(I) = \lambda \Delta^2(I)|A(I)|/(\pi_{-p}^{\bar{\sigma}^t}(I)t)$. Put $u_p^{\bar{\sigma}^t}(I, a)$ into (8), we get

$$u_p^{\bar{\sigma}^t}(I, a) = \mathbb{E}_{h \sim I, h' \sim \text{Succ}_{\bar{\sigma}^t}(h, a)} \left\{ \mathbb{1}_{h' \in Z} u_p(h') + \mathbb{1}_{h' \notin Z} u_p^{\bar{\sigma}^t}(I(h')) \right\}. \quad (23)$$

A proof for the equation is provided in Appendix C. Equation (23) implies that if we *sample* the payoffs and the RSVs of the earliest reachable histories starting from (I, a) , then, the expectation of the sampled values is precisely the RSV of action a at infoset I . According to this equation, a bootstrap method similar to Q-learning [36] is derived.

A. Bootstrap Learning for RSVs

According to Equation (23), we propose ReCFR with bootstrapping (ReCFR-B).

Definition 1. ReCFR-B is an ReCFR that learns the RSVs using a bootstrap method: at each iteration t of ReCFR-B, initialize $u_p^{\bar{\sigma}^t, 1}(I, a)$ with an arbitrary value for each $p \in P$, $I \in \mathcal{I}_p$, and $a \in A(I)$. Let each $p \in P$ play K games with the opponent who uses strategy $\bar{\sigma}_{-p}^t$. During the play, for each transition (h, a, h') encountered in game $1 \leq k \leq K$, update $u_p^{\bar{\sigma}^t, k}(I(h), a)$ according to

$$u_p^{\bar{\sigma}^t, k+1}(I(h), a) = u_p^{\bar{\sigma}^t, k}(I(h), a) + \alpha_k \left(\mathbb{1}_{h' \in Z} u_p(h') + \mathbb{1}_{h' \notin Z} \gamma u_p^{\bar{\sigma}^t, k}(I(h')) - u_p^{\bar{\sigma}^t, k}(I(h), a) \right),$$

where $u_p^{\bar{\sigma}^t, k}(I)$ fulfills

$$\sum_{a \in A(I)} \left(u_p^{\bar{\sigma}^t, k}(I, a) - u_p^{\sigma^t}(I) \right)_+^2 = \beta_p^t(I).$$

Theorem 3. At each iteration t of ReCFR-B, if every terminal history is visited with a non-zero probability in each game, $0 < \gamma \leq 1$, $\sum_k \alpha_k = \infty$, and $\sum_k \alpha_k^2 < \infty$, then, $u_p^{\bar{\sigma}^t, K}(I, a)$ converges to $u_p^{\bar{\sigma}^t}(I, a)$ w.p.1 for every (I, a) as $K \rightarrow \infty$.

The proof is provided in Appendix C. In practice, we set α_k to a constant learning rate and set γ to 1 as NFSP did [28]. According to the definition, the update of $u_p^{\bar{\sigma}^t, k}(I, a)$ does not depend on how infoset I is reached and how the player p acts in descendants. So the learning is **off-policy**. ReCFR-B is similar to (batch-)OSCFR [33], as both of them update the values of infosets by sampling trajectories. OSCFR

is usually considered model-free [20], [29].⁶ So, if $\beta_p^t(I)$ for each infoset is chosen without referring to any private information, e.g., $\sigma_c(h, a)$ and $\Delta(I)$, about the transition model, we can consider ReCFR-B **model-free**. In practice, we use $\beta_p^t(I) = \lambda \Delta^2(I)|A(I)|/(\pi_{-p}^{\bar{\sigma}^t}(I)t)$ as it gives a good exploitability bound.

B. Neural ReCFR with Bootstrapping

For each player, an RSV network $\mathcal{R}(\theta_p^t)$ with parameters θ_p^t is used to approximate $u_p^{\bar{\sigma}^t}(I, a)$ at all infosets. Also, a neural network $\Pi(\phi_p^t)$ with parameters ϕ_p^t is used to approximate the average strategy $\bar{\sigma}_p^t$. At iteration t , the next iteration strategy σ_p^{t+1} is computed according to the output of the RSV network,

$$\sigma_p^{t+1}(I, a) = \text{RM} \left(\mathcal{R}(I, a | \theta_p^t) - u_p^{\bar{\sigma}^t}(I) \right), \quad (24)$$

where $u_p^{\bar{\sigma}^t}(I)$ is the RSV of infoset I , and it is chosen to fulfill (22) with $u_p^{\bar{\sigma}^t}(I, a) \leftarrow \mathcal{R}(I, a | \theta_p^t)$. Since RM is *scale-invariant*, this equation is equivalent to (17) if the approximation is accurate. We use anticipatory dynamics [28] to estimate the average strategy $\bar{\sigma}^t$. Specifically, the strategy for each player at iteration t is $\hat{\sigma}_p = (1 - \eta)\Pi(\phi_p^{t-1}) + \eta\sigma_p^t$, where η is a hyper-parameter. We set $\eta = 0.1$ as NFSP did. In practice, multiple trajectories are sampled. So, $\Pi(\phi_p^{t-1})$ is selected with probability $1 - \eta$ for playing, while σ_p^t is selected with probability η . The pseudocode of Neural ReCFR-B is given in Algorithm 2. At each iteration and for each player p , K games are played using the anticipatory strategy $\hat{\sigma}$. During the play, the transitions are collected into a temporary data set $\mathcal{D}_p^{\mathcal{R}}$, and the behavior tuples are stored in a reservoir buffer [28] \mathcal{M}_p^{Π} . The RSV network for player p is trained by minimizing:

$$\mathcal{L}(\theta_p^t) = \mathbb{E}_{(h, a, h', \pi_{-p}^{\bar{\sigma}^t}(I(h))) \sim \mathcal{D}_p^{\mathcal{R}}} \left\{ \left(\mathbb{1}_{h' \in Z} u_p(h') + \mathbb{1}_{h' \notin Z} u_p^{\bar{\sigma}^t}(I(h')) - \mathcal{R}(I(h), a | \theta_p^t) \right)^2 \right\}, \quad (25)$$

where θ_p^t is the trainable parameters initialized with θ_p^{t-1} , and $u_p^{\bar{\sigma}^t}(I(h'))$ is chosen to fulfill (22) with $u_p^{\bar{\sigma}^t}(I(h'), a) \leftarrow \mathcal{R}(I(h'), a | \theta_p^t)$. The average network for player p is trained by minimizing the cross-entropy loss:

$$\mathcal{L}(\phi_p^t) = \mathbb{E}_{(h, a) \sim \mathcal{M}_p^{\Pi}} \left\{ -\log \Pi(I(h), a | \phi_p^t) \right\}, \quad (26)$$

where ϕ_p^t is the trainable parameters initialized with ϕ_p^{t-1} .

The algorithm requires each player $p \in P$ plays with an opponent who uses the average strategy $\bar{\sigma}_{-p}^t$. However, the strategy for player p is not specified. We propose **asymmetric learning**, in which player p uses a mixed strategy of the uniform random strategy and strategy σ_p^t , i.e., $\hat{\sigma}_p = (1 - \eta)\text{Uniform} + \eta\sigma_p^t$. Accordingly, the method that both players use the anticipatory strategy is called **symmetric learning**. We use asymmetric learning in our experiments by default.

⁶The players are regarded as a whole. However, the chance player's strategy is a part of the transition model.

Algorithm 2 Neural ReCFR with bootstrapping

```

1: Input: maximum iterations  $T$ .
2: Initialize RSV network parameters  $\theta_p^0$  and policy network parameters  $\phi_p^0$  for each player.
3: Initialize strategy memory buffer  $\mathcal{M}_p^\Pi$  for each player.
4: for iteration  $t = 1$  to  $T$  do
5:   for player  $p \in P$  do
6:     Initialize a temporary data set  $\mathcal{D}_p^{\mathcal{R}}$ .
7:     for game  $k = 1$  to  $K$  do
8:       Sample a strategy  $\hat{\sigma}_i$  for each player  $i \in P$ ,  $\hat{\sigma}_i = \begin{cases} \sigma_i^t, & \text{with probability } \eta, \\ \Pi(\phi_i^{t-1}), & \text{with probability } 1 - \eta. \end{cases}$ 
9:        $\mathcal{T}_p^k \leftarrow \text{Play}(\hat{\sigma}_p, \hat{\sigma}_{-p})$ . ▷ Sample a trajectory  $\mathcal{T}_p^k$  using  $\hat{\sigma}$ .
10:      Collect all transitions  $(h, a, h', \pi_{-p}^{\hat{\sigma}^t}(I(h))) \in \mathcal{T}_p^k$  to data set  $\mathcal{D}_p^{\mathcal{R}}$ .
11:      if player  $p$  follows the current strategy  $\sigma^t$  then
12:        Store all behavior tuples  $(h, a) \in \mathcal{T}_p^k$  in strategy memory buffer  $\mathcal{M}_p^\Pi$ .
13:      Train  $\theta_p^t$  on loss (25) and train  $\phi_p^t$  on loss (26).
```

C. Discussion

Approximating the average strategy. In Neural ReCFR-B, approximating cumulative regrets is avoided, but the average strategy is still approximated. Single Deep CFR [27] proposed remembering all past regret networks to avoid doing that. This method should also be compatible with Neural ReCFR-B. Nevertheless, it should be easier to train the average networks since the average strategy changes more slowly than the cumulative regrets.

Relationship with DNCFR. Both Neural ReCFR-B and DNCFR use bootstrap learning to reduce the variance in training targets. However, a critical difference is that Neural ReCFR-B bootstraps on the RSVs of the earliest reachable infosets, while DNCFR bootstraps on the last approximations. Since we only consider depth-limited games, the approximation error in Neural ReCFR-B is *bounded* at every iteration, while the error in DNCFR is *accumulated*.

Relationship with NFSP. Since ReCFR is equivalent to XFP according to Proposition 1 when $\lambda_p^t(I) = 0$ at all infosets, Neural ReCFR-B can also be regarded as NFSP in this case. Note that the learning algorithm, i.e., DQN [3], in NFSP has two essential components: 1) a value memory buffer, which may improve sample efficiency. 2) a target network, which may stabilize the training. We will test Neural ReCFR-B with these two components in section VII.

VII. EXPERIMENTAL SETUP AND RESULTS

We first test ReCFR on Leduc Poker [37] to show the convergence properties of the algorithm. Leduc poker is a small size game with two rounds of betting. Then, Neural ReCFR-B is tested on heads-up flop hold'em poker (FHP) [18] and heads-up limit Texas hold'em (HULH).⁷ FHP is a medium-size game with over 10^{12} nodes and 10^9 infosets, while HULH is a large-size game with over 10^{17} nodes and 10^{14} infosets. More details about the games are given in Appendix D. The neural network architecture is the same as in [18]. It is a seven-layer fully connected neural network.

We train the neural networks using Adam optimizer [38], with a batch size of 128, a learning rate of 0.001, and gradient norm clipping to 1. At every iteration, 1000 plays are performed to collect samples. The RSV networks are trained for two epochs (approximately 32 SGD steps with a batch size of 128 on FHP), while the average networks are trained for 16 SGD steps. The strategy memory size is set to 10 million. For HULH, we increase the batch size to 6,400, the number of SGD steps to 64, and the number of plays to 100,000. The strategy memory size is also increased to 40 million. We also test Neural ReCFR-B with two RSV memory buffers (one for each player). In this setting, the RSV memory sizes for FHP and HULH are 1 million and 4 million, respectively, and the RSV networks are trained for 32 and 64 SGD steps, respectively.

Neural ReCFR-B is compared with two model-based neural CFRs: Deep CFR and DNCFR; and two model-free algorithms: DREAM and NFSP, on FHP. A model-free variant of Deep CFR with Outcome Sampling [33] (Deep OSCFR) is also included. All the algorithms are implemented based on OpenSpiel [39]. We implement Deep CFR and DREAM with the hyper-parameters given in [18] and [20], respectively. The hyper-parameters for DNCFR and NFSP are determined through a set of experiments. A head-to-head comparison is performed between Neural ReCFR-B and Deep CFR on HULH. All the settings are given in Appendix E.

Performance is measured using exploitability in terms of milli big blinds per game (mbb/g). For a specific exploitability value, the numbers of nodes touched and the numbers of samples consumed (= #SGD steps \times batch size) of different algorithms are compared, respectively. Algorithms touching fewer nodes are more sample efficient, while algorithms consuming fewer samples are more training efficient.

A. Experimental Results of ReCFR

In this subsection, we test ReCFR on Leduc poker [37]. ReCFR is compared with vanilla CFR, as the latter can be regarded as a special case of ReCFR. We first test ReCFR with a constant λ in $[0, 1]$. As shown on the left side in Figure 1,

⁷Source code: https://github.com/Liuweiming/Neural_ReCFR_B

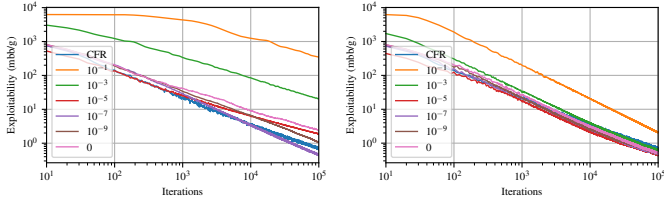


Fig. 1: Exploitability curves of vanilla CFR and ReCFR on Leduc poker. **Left:** ReCFR with constant λ s shown in the legend. **Right:** ReCFR with adaptive λ s initialized with the values in the legend.

ReCFR with any λ converges. It seems ReCFR with $\lambda = 10^{-7}$ is the fastest, even faster than vanilla CFR. However, if λ is reduced to 0, ReCFR will degenerate to XFP, and it converges slower than CFR. We also test ReCFR with an adaptive λ , as shown on the right side in Figure 1. As we can see, ReCFR with different initial λ (even 0) converges as fast as CFR, except the one with an initial $\lambda = 1$. So, with the adaptive method, it is easier to choose the hyper-parameter for ReCFR, as long as the initial value is small enough. In Appendix F, the curves of the adaptive λ in ReCFR with different initial values are given.

B. Experimental Results of Neural ReCFR-B

The results of all the algorithms on FHP are given in Figure 2. As we can see, Neural ReCFR-B achieves an exploitability of 50.5 mbb/g after touching 1.0×10^{10} nodes, while Deep CFR achieves 47.0 mbb/g after touching 6.0×10^8 nodes, but the value increases to 73.0 mbb/g after touching 1.3×10^9 nodes. Neural ReCFR-B is worse than Deep CFR in sample efficiency. It is reasonable because model-free algorithms are less sample efficient by nature. However, our algorithm is faster than the three model-free algorithms: Deep OSCFR, DREAM, and NFSP. More importantly, our algorithm is the most training efficient, by more than **25** times faster than Deep CFR to reach the exploitability of 100 mbb/g, while Deep OSCFR and DREAM never reach this value and both DNCFR and NFSP are stuck at this value. Recall that DNCFR also uses bootstrap learning. However, it has a higher exploitability than Neural ReCFR-B. On the other hand, NFSP is training efficient, but it achieves an exploitability much higher than Neural ReCFR-B and Deep CFR. Note that NFSP can be regarded as a special case of Neural ReCFR-B to some extent, and the latter may degenerate to NFSP when the gap between $u_p^{\sigma^t}(I)$ and $\max_a u_p^{\sigma^t}(I, a)$ at every infoset is too small to be learned by neural networks. As for DREAM, it performs better than Deep OSCFR but worse than Neural ReCFR-B and Deep CFR. This suggests that DREAM is still suffering from high variance training targets.

We also tried to improve the training efficiency of Deep CFR by only reducing the number of SGD steps for training the regret networks. The results in Appendix F show that Deep CFRs with fewer SGD steps converge earlier to higher exploitability. On the other hand, as shown in Figure 3, Neural ReCFR-B converges to a similar exploitability even the RSV networks are trained for one epoch per iteration. To further investigate the robustness of Neural ReCFR-B, we test it with different settings on FHP. In Figure 3, the results of

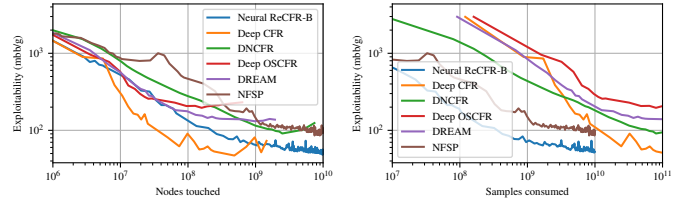


Fig. 2: Exploitability curves of different algorithms on FHP. The x-axes represent the number of nodes touched and the number of samples consumed (= #SGD steps \times batch size), respectively.

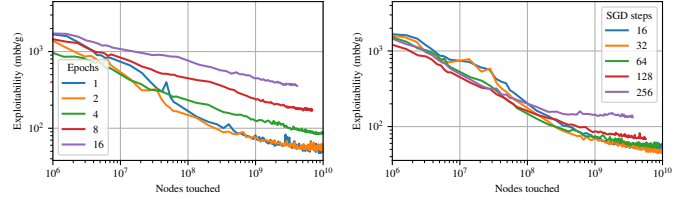


Fig. 3: Exploitability curves of Neural ReCFR-B with different numbers of epochs / SGD steps on FHP. **Left:** the default Neural ReCFR-B. **Right:** Neural ReCFR-B with RSV memory buffers.

Neural ReCFR-B with different numbers of epochs for training the RSV networks are given. As we can see, training the RSV networks for more than two epochs can reduce the convergence speed significantly. A possible reason is that the neural networks are overfitting. We conjecture that using memory buffers can partially solve this problem. As shown on the right side in Figure 3, the problem is alleviated. Overfitting is also observed in Deep CFR [18], and it is solved by training the regret networks from scratch at every iteration. In Neural ReCFR-B, reducing the number of epochs to alleviate the problem is more appealing. In Figure 4, the results of Neural ReCFR-B with different numbers of plays per iteration are given. It is shown that Neural ReCFR-B is insensitive to the number of plays in sample efficiency. However, performing more plays (collecting more training data) per iteration can increase the training efficiency. This is because of the variance in approximated RSVs is also affected by the data set size.

On the left side in Figure 5, we test the scalability of the algorithm by increasing the number of plays per iteration and the batch size simultaneously. As we can see, the algorithm can scale down to 500 plays and up to 50,000 plays per iteration. Performing more plays per iteration and using a larger batch size may reduce the sample efficiency. However, it is easier to parallelize and thus may reduce the training time dramatically. In conclusion, Neural ReCFR-B is robust to different hyper-parameters.

On the right side in Figure 5, we test Neural ReCFR-B with additional components. The default is Neural ReCFR-B with asymmetric learning, without memory buffers and target networks. As we can see, Neural ReCFR-B with target networks performs identically to the default one. On the other hand, Neural ReCFR-B with RSV memory buffers seems to have a slightly lower exploitability. This result suggests that using RSV memory buffers in Neural ReCFR-B may increase the sample efficiency. However, Neural ReCFR-B with symmetric learning is worse than the default one. The reason might be that symmetric learning can not guarantee to

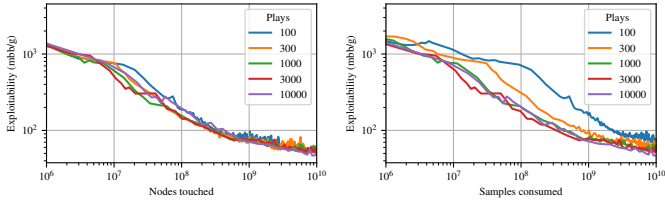


Fig. 4: Exploitability curves of Neural ReCFR-B with different numbers of plays on FHP. The x-axes represent the number of nodes touched and the number of samples consumed, respectively.

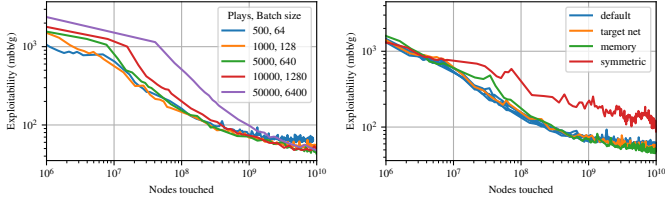


Fig. 5: **Left:** Exploitability curves of Neural ReCFR-B with different numbers of plays and batch sizes. **Right:** Exploitability curves of Neural ReCFR-B with additional components. We also plot the exploitability curves of five replicates of Neural ReCFR-B (the blue lines) with different random seeds.

visit every infoset with a probability greater than zero. Besides, we plot five replicates of Neural ReCFR-B with different random seeds (so the neural networks are also initialized differently). As we can see, the results are consistent.

Finally, we evaluate Neural ReCFR-B on HULH. Since using RSV memory buffers may increase sample efficiency and avoid overfitting, we use RSV memory buffers in this experiment. We train Neural ReCFR-B and Deep CFR separately for 20 days on one GPU and 10 CPU cores. On the left side of Figure 6, we compare the lower bounds on the exploitability of Neural ReCFR-B and Deep CFR. The lower bounds are estimated using a Local Best Response [40] algorithm. The comparison is performed day by day. So it should reflect both sample efficiency and training efficiency. As we can see, Neural ReCFR-B achieves a lower bound around 600 mbb/g only after one day’s training, while Deep CFR reaches this value on the ninth day. After 20 days’ training, Deep CFR and Neural ReCFR-B have similar lower bounds on exploitability, about 300 mbb/g. Besides, Figure 7 shows that Neural ReCFR-B is more training efficient, but Deep CFR has a better sample efficiency. Note that the sample efficiency of Neural ReCFR-B may be increased when the batch size is reduced at the cost of increasing training time, according to the scalability experiment on FHP. On the right side of Figure 6, we show the head-to-head performance between Neural ReCFR-B and Deep CFR. As we can see, Neural ReCFR-B beats Deep CFR by more than 150 mbb/g. It looks like that both algorithms have not converged yet, but it is safe to conclude that Neural ReCFR-B is more training efficient and time-efficient.

VIII. CONCLUSION

This paper proposes a new CFR algorithm, ReCFR, and a model-free Neural CFR algorithm, Neural ReCFR-B. In these two algorithms, cumulative regrets are replaced by RSVs proposed in [22]. After revisiting RSVs and Warm CFR, we

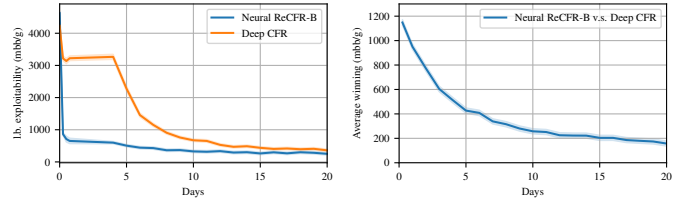


Fig. 6: Comparison of Neural ReCFR-B and Deep CFR on HULH. The x-axis indicates the days for training. **Left:** Lower bound on exploitability. **Right:** Head-to-head performance. The y-axis represents the average winning of Neural ReCFR-B (simulated by 10^6 plays). All the values are presented with 95% confidence interval.

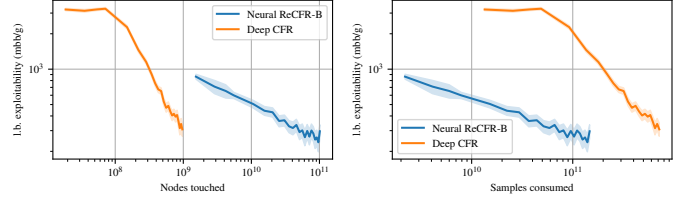


Fig. 7: Curves of lower bounds on exploitability of Neural ReCFR-B and Deep CFR on HULH. The x-axes represent the number of nodes touched and the number of samples consumed, respectively.

prove that ReCFR can converge to a Nash equilibrium at a rate of $O(1/\sqrt{T})$. Thanks to the recursive and non-cumulative properties of the RSVs, when bootstrap learning is used, the variance in training targets in Neural ReCFR-B should be low. The experimental results show that Neural ReCFR-B achieves competitive performance to the state-of-the-art neural CFR algorithms with higher training efficiency.

According to the theoretical results, it is promising to transform other CFR algorithms, e.g., CFR+ [26] and PCFR [41], to new algorithms similar to ReCFR, and develop new neural CFR algorithms based on them. Also, combining Neural ReCFR-B with variance-reduction techniques [30], [42] or other improvements [27], [20], [29] may produce more efficient algorithms. Furthermore, there are some other equilibrium-finding algorithms that have fast convergence, e.g., Excessive Gap Technique (EGT) [11] and Optimistic Follow the Regularized Lead (OFTRL) [43]. And EGT does not even use any cumulative variables. However, it is unclear whether they are compatible with sampling and function approximation. So, more research is required in this direction. Besides, it has been shown in [44] that RM is equivalent to Follow the Regularized Lead (FTRL) [45]. In this paper, it is shown that ReCFR is a generalization of CFR and XFP. So, there might be some strong connections between CFR, XFP, and FTRL worthy of further study.

With the model-free Neural ReCFR-B, we may also apply CFR algorithms to a broader range of IIGs, for example, non-stationary games, non-zero-sum games, multi-player games, and even video games.

REFERENCES

- [1] G. Tesauro, “Temporal difference learning and td-gammon,” *Commun. ACM*, vol. 38, no. 3, pp. 58–68, 1995.
- [2] R. Coulom, “Efficient selectivity and backup operators in monte-carlo tree search,” in *Computers and Games, 5th International Conference, CG 2006, Turin, Italy, May 29-31, 2006. Revised Papers*, ser. Lecture Notes in Computer Science, vol. 4630, 2006, pp. 72–83.
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [4] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. P. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, “Mastering the game of go without human knowledge,” *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [5] P. I. Cowling, E. J. Powley, and D. Whitehouse, “Information set monte carlo tree search,” *IEEE Trans. Comput. Intell. AI Games*, vol. 4, no. 2, pp. 120–143, 2012.
- [6] N. Brown and T. Sandholm, “Safe and nested subgame solving for imperfect-information games,” in *Advances in Neural Information Processing Systems 30*, 2017, pp. 689–699.
- [7] S. Srinivasan, M. Lanctot, V. F. Zambaldi, J. Pérolat, K. Tuyls, R. Munos, and M. Bowling, “Actor-critic policy optimization in partially observable multiagent environments,” in *Advances in Neural Information Processing Systems*, 2018, pp. 3426–3439.
- [8] M. Zinkevich, M. Johanson, M. H. Bowling, and C. Piccione, “Regret minimization in games with incomplete information,” in *Advances in Neural Information Processing Systems 20*, 2007, pp. 1729–1736.
- [9] J. Heinrich, M. Lanctot, and D. Silver, “Fictitious self-play in extensive-form games,” in *International Conference on Machine Learning*, vol. 37, 2015, pp. 805–813.
- [10] B. Bosanský, C. Kiekintveld, V. Lisý, J. Cermak, and M. Pechoucek, “Double-oracle algorithm for computing an exact nash equilibrium in zero-sum extensive-form games,” in *International conference on Autonomous Agents and Multi-Agent Systems*, 2013, pp. 335–342.
- [11] S. Hoda, A. Gilpin, J. Peña, and T. Sandholm, “Smoothing techniques for computing nash equilibria of sequential games,” *Mathematics of Operations Research*, vol. 35, no. 2, pp. 494–512, 2010.
- [12] M. Moravčík, M. Schmid, N. Burch, V. Lisý, D. Morrill, N. Bard, T. Davis, K. Waugh, M. Johanson, and M. Bowling, “Deepstack: Expert-level artificial intelligence in heads-up no-limit poker,” *Science*, vol. 356, no. 6337, pp. 508–513, 2017.
- [13] N. Brown and T. Sandholm, “Superhuman AI for heads-up no-limit poker: Libratus beats top professionals,” *Science*, vol. 359, no. 6374, pp. 418–424, 2018.
- [14] N. Brown and T. Sandholm, “Superhuman AI for multiplayer poker,” *Science*, vol. 365, no. 6456, pp. 885–890, 2019.
- [15] M. Johanson, N. Burch, R. A. Valenzano, and M. Bowling, “Evaluating state-space abstractions in extensive-form games,” in *International conference on Autonomous Agents and Multi-Agent Systems*, 2013, pp. 271–278.
- [16] K. Waugh, D. Morrill, J. A. Bagnell, and M. H. Bowling, “Solving games with functional regret estimation,” in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015, pp. 2138–2145.
- [17] D. R. Morrill, “Using regret estimation to solve games compactly,” *Master’s thesis, University of Alberta*, 2016.
- [18] N. Brown, A. Lerer, S. Gross, and T. Sandholm, “Deep counterfactual regret minimization,” in *International Conference on Machine Learning*, vol. 97, 2019, pp. 793–802.
- [19] H. Li, K. Hu, S. Zhang, Y. Qi, and L. Song, “Double neural counterfactual regret minimization,” in *International Conference on Learning Representations*, 2020.
- [20] E. Steinberger, A. Lerer, and N. Brown, “DREAM: deep regret minimization with advantage baselines and model-free learning,” *CoRR*, vol. abs/2006.10410, 2020.
- [21] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2nd Edition*, ser. Springer Series in Statistics. Springer, 2009.
- [22] N. Brown and T. Sandholm, “Strategy-based warm starting for regret minimization in games,” in *AAAI Conference on Artificial Intelligence*, 2016, pp. 432–438.
- [23] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [24] N. Burch, M. Johanson, and M. Bowling, “Solving imperfect information games using decomposition,” in *AAAI Conference on Artificial Intelligence*, 2014, p. 602–608.
- [25] N. Brown, T. Sandholm, and B. Amos, “Depth-limited solving for imperfect-information games,” in *Advances in Neural Information Processing Systems 31*, 2018, pp. 7663–7674.
- [26] O. Tammelin, “Solving large imperfect information games using CFR+,” *CoRR*, vol. abs/1407.5042, 2014.
- [27] E. Steinberger, “Single deep counterfactual regret minimization,” *CoRR*, vol. abs/1901.07621, 2019.
- [28] J. Heinrich and D. Silver, “Deep reinforcement learning from self-play in imperfect-information games,” *CoRR*, vol. abs/1603.01121, 2016.
- [29] A. Gruslys, M. Lanctot, R. Munos, F. Timbers, M. Schmid, J. Pérolat, D. Morrill, V. F. Zambaldi, J. Lespiau, J. Schultz, M. G. Azar, M. Bowling, and K. Tuyls, “The advantage regret-matching actor-critic,” *CoRR*, vol. abs/2008.12234, 2020.
- [30] M. Schmid, N. Burch, M. Lanctot, M. Moravcik, R. Kadlec, and M. Bowling, “Variance reduction in monte carlo counterfactual regret minimization (VR-MCCFR) for extensive form games using baselines,” in *AAAI Conference on Artificial Intelligence*, 2019, pp. 2157–2164.
- [31] S. Srinivasan, M. Lanctot, V. F. Zambaldi, J. Pérolat, K. Tuyls, R. Munos, and M. Bowling, “Actor-critic policy optimization in partially observable multiagent environments,” in *Advances in Neural Information Processing Systems*, 2018, pp. 3426–3439.
- [32] E. Lockhart, M. Lanctot, J. Pérolat, J. Lespiau, D. Morrill, F. Timbers, and K. Tuyls, “Computing approximate equilibria in sequential adversarial games by exploitability descent,” in *International Joint Conference on Artificial Intelligence*, 2019, pp. 464–470.
- [33] M. Lanctot, K. Waugh, M. Zinkevich, and M. H. Bowling, “Monte carlo sampling for regret minimization in extensive games,” in *Advances in Neural Information Processing Systems*, 2009, pp. 1078–1086.
- [34] M. Johanson, N. Bard, M. Lanctot, R. G. Gibson, and M. Bowling, “Efficient nash equilibrium approximation through monte carlo counterfactual regret minimization,” in *International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2012, Valencia, Spain, June 4-8, 2012 (3 Volumes)*, 2012, pp. 837–846.
- [35] R. G. Gibson, M. Lanctot, N. Burch, D. Szafron, and M. Bowling, “Generalized sampling and variance in counterfactual regret minimization,” in *AAAI Conference on Artificial Intelligence*, 2012.
- [36] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, no. 3-4, pp. 279–292, 1992.
- [37] F. Southey, M. Bowling, B. Larson, C. Piccione, N. Burch, D. Billings, and C. Rayner, “Bayes’ bluff: opponent modelling in poker,” in *Conference on Uncertainty in Artificial Intelligence*, 2005, pp. 550–558.
- [38] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representation*, 2015.
- [39] M. Lanctot, E. Lockhart, J. Lespiau, V. F. Zambaldi, S. Upadhyay, J. Pérolat, S. Srinivasan, F. Timbers, K. Tuyls, S. Omidshafiei, D. Hennes, D. Morrill, P. Muller, T. Ewalds, R. Faulkner, J. Kramár, B. D. Vyllder, B. Saeta, J. Bradbury, D. Ding, S. Borgeaud, M. Lai, J. Schrittwieser, T. W. Anthony, E. Hughes, I. Danihelka, and J. Ryan-Davis, “OpenSpiel: A framework for reinforcement learning in games,” *CoRR*, vol. abs/1908.09453, 2019.
- [40] V. Lisý and M. Bowling, “Equilibrium approximation quality of current no-limit poker bots,” in *The Workshops of the AAAI Conference on Artificial Intelligence*, vol. WS-17, 2017.
- [41] G. Farina, C. Kroer, and T. Sandholm, “Faster game solving via predictive blackwell approachability: Connecting regret matching and mirror descent,” in *AAAI Conference on Artificial Intelligence*, 2021, pp. 5363–5371.
- [42] T. Davis, M. Schmid, and M. Bowling, “Low-variance and zero-variance baselines for extensive-form games,” in *International Conference on Machine Learning*, vol. 119, 2020, pp. 2392–2401.
- [43] G. Farina, C. Kroer, and T. Sandholm, “Optimistic regret minimization for extensive-form games via dilated distance-generating functions,” in *Advances in Neural Information Processing Systems*, 2019, pp. 5222–5232.
- [44] K. Waugh and J. A. Bagnell, “A unified view of large-scale zero-sum equilibrium computation,” in *The Workshops of the AAAI Conference on Artificial Intelligence*, ser. AAAI Technical Report, S. Ganzfried, Ed., vol. WS-15-07. AAAI Press, 2015.
- [45] J. D. Abernethy, E. Hazan, and A. Rakhlin, “Competing in the dark: An efficient algorithm for bandit linear optimization,” in *Conference on Learning Theory*, 2008, pp. 263–274.
- [46] L. Condat, “Fast projection onto the simplex and the l_1 ball,” *Math. Program.*, vol. 158, no. 1-2, pp. 575–585, 2016.

- [47] D. S. Leslie and E. J. Collins, "Generalised weakened fictitious play," *Games Econ. Behav.*, vol. 56, no. 2, pp. 285–298, 2006.
- [48] F. Orabona, "A modern introduction to online learning," *arXiv preprint arXiv:1912.13213*, 2019.
- [49] T. S. Jaakkola, M. I. Jordan, and S. P. Singh, "On the convergence of stochastic iterative dynamic programming algorithms," *Neural Comput.*, vol. 6, no. 6, pp. 1185–1201, 1994.
- [50] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. A. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zhang, "Tensorflow: A system for large-scale machine learning," *CoRR*, vol. abs/1605.08695, 2016.

IX. PROOFS AND MORE DETAILS FOR CFR AND RSVs

A. Proof for Equation 2

Proof. First, for any strategy σ , a history of player p , a descendant $h' \in Succ_p(h \cdot a)$, and any terminal history $z \in Z$ reachable from h and h' , according to the definition of *reach*, we have

$$\begin{aligned} \pi^\sigma(h \cdot a, z) &= \pi_{-p}^\sigma(h \cdot a, z) \times \pi_p^\sigma(h \cdot a, z) \\ &= \pi_{-p}^\sigma(h \cdot a, h') \pi_{-p}^\sigma(h', z) \times \pi_p^\sigma(h', z) \quad (\text{A.27}) \\ &= \pi_{-p}^\sigma(h \cdot a, h') \pi^\sigma(h', z). \end{aligned}$$

Note that $\pi_p^\sigma(h \cdot a, z) = \pi_p^\sigma(h', z)$ according to the definition of *Succ_p*. Therefore,

$$\begin{aligned} \pi_{-p}^\sigma(h) \pi^\sigma(h \cdot a, z) &= \pi_{-p}^\sigma(h) \pi_{-p}^\sigma(h \cdot a, h') \pi^\sigma(h', z) \\ &= \pi_{-p}^\sigma(h') \pi^\sigma(h', z). \quad (\text{A.28}) \end{aligned}$$

So, for any strategy σ , $I \in \mathcal{I}_p$ and $a \in A(I)$,

$$\begin{aligned} &v_p^\sigma(I, a) \\ &= \sum_{h \in I} \sum_{z \in Z: h \sqsubset z} \pi_{-p}^\sigma(h) \pi^\sigma(h \cdot a, z) u_p(z) \\ &= \sum_{h \in I} \sum_{z \in Z: z \in Succ_p(h \cdot a)} \pi_{-p}^\sigma(z) \pi^\sigma(z, z) u_p(z) + \\ &\quad \sum_{h \in I} \sum_{h' \notin Z: h' \in Succ_p(h \cdot a)} \sum_{z \in Z: h' \sqsubset z} \pi_{-p}^\sigma(h') \pi^\sigma(h', z) u_p(z) \\ &= \sum_{h \in I} \sum_{z \in Z: z \in Succ_p(h \cdot a)} \pi_{-p}^\sigma(z) u_p(z) + \\ &\quad \sum_{I' \in Succ_p(I, a)} \sum_{h' \in I'} \sum_{z \in Z: h' \sqsubset z} \pi_{-p}^\sigma(h') \pi^\sigma(h', z) u_p(z) \\ &= \sum_{h \in I} \sum_{z \in Z: z \in Succ_p(h \cdot a)} \pi_{-p}^\sigma(z) u_p(z) + \\ &\quad \sum_{I' \in Succ_p(I, a)} v_p^\sigma(I'). \quad (\text{A.29}) \end{aligned}$$

Note that we only consider perfect-recall games. The third equality holds because $\{h' | h' \notin Z, h' \in Succ_p(h \cdot a), h \in I\} = \{h' | h' \in I', I' \in Succ_p(I, a)\}$. Then, the fourth equality holds according to the definition of counterfactual values. \square

B. A Method for Computing the RSVs

When computing the RSVs, we need to solve $v_p^{\prime\sigma}(I) \in \mathbb{R}$ in constraint

$$\sum_{a \in A(I)} (v_p^{\prime\sigma}(I, a) - v_p^{\prime\sigma}(I))_+^2 \leq \lambda(I), \quad (\text{A.30})$$

where $\lambda(I) \in (0, \infty)$ is a parameter. There could be multiple solutions that satisfy the constraint. Define $f : v_p^{\prime\sigma}(I) \mapsto \sum_{a \in A(I)} (v_p^{\prime\sigma}(I, a) - v_p^{\prime\sigma}(I))_+^2$. The constraint can be rewritten as $f(v_p^{\prime\sigma}(I)) \leq \lambda(I)$. To solve the inequation, we need to first solve $v_p^{\prime\sigma}(I)$ in equation

$$f(v_p^{\prime\sigma}(I)) = \lambda(I). \quad (\text{A.31})$$

We plotted $f(x)$ in Figure 8, in the case that $v_p^{\prime\sigma}(I, \cdot) = [-0.7, 0, 1]$. As we can see, function $f(x)$ in range $(-\infty, \max_a v_p^{\prime\sigma}(I, a)]$ is strongly convex and strictly decreasing. In other words, function $f : (-\infty, \max_a v_p^{\prime\sigma}(I, a)] \rightarrow [0, \infty)$ is a bijection. Therefore, the solution for (A.31) exists and is unique when $\lambda(I) > 0$.

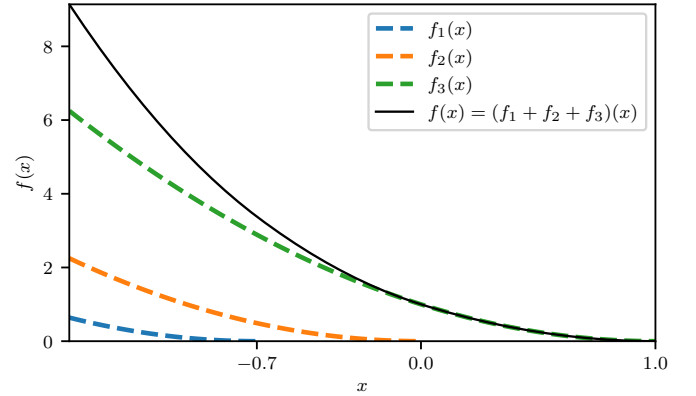


Fig. 8: Curves of $f(x)$ when $v_p^{\prime\sigma}(I, \cdot) = [-0.7, 0, 1]$.

To solve (A.31), we can first sort $v_p^{\prime\sigma}(I, \cdot)$ such that $v_p^{\prime\sigma}(I, a_1) \leq \dots \leq v_p^{\prime\sigma}(I, a_{|A(I)|})$. Then we can try to solve $v_p^{\prime\sigma}(I)$ in each interval: $(\infty, v_p^{\prime\sigma}(I, a_1)]$, \dots , $(v_p^{\prime\sigma}(I, a_{|A(I)|-1}), v_p^{\prime\sigma}(I, a_{|A(I)|})]$. In each interval, the hard thresholding operator $(\cdot)_+$ can be removed, i.e., (A.31) is equivalent to a conventional quadratic equation.

Since the complexity for sorting $v_p^{\prime\sigma}(I, \cdot)$ and iteratively solving the quadratic equations is $O(|A(I)|^2)$, The complexity for solving (A.30) or (A.31) is $O(|A(I)|^2)$. A similar problem can be found in [46], where we may find an algorithm to solve (A.30) or (A.31) more efficiently.

C. Proof for Lemma 2

Proof. The right hand side of the equation is

$$\begin{aligned} &\sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) (v_p^\sigma(I, a) - v_p^\sigma(I)) \sigma_p'(I, a) \\ &= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) v_p^\sigma(I, a) \sigma_p'(I, a) - \\ &\quad \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) v_p^\sigma(I) \quad (\text{A.32}) \end{aligned}$$

Note that $\pi_p^{\sigma'}(I) \sigma_p'(I, a) = \pi_p^{\sigma'}(h) \sigma_p'(h, a) = \pi_p^{\sigma'}(h') = \pi_p^{\sigma'}(I')$ for any $h \in I, h' \in Succ_p(h \cdot a)$ and $I' \in Succ_p(I, a)$. So, for the first term on the right side in the above equation,

according to the recursive definition of counterfactual values (Equation (2) in the paper),

$$\begin{aligned}
& \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) v_p^\sigma(I, a) \sigma'_p(I, a) \\
&= \sum_{I \in \mathcal{I}_p} \sum_{h \in I} \sum_{a \in A(I)} \sum_{z \in Z: z \in \text{Succ}_p(h, a)} \pi_p^{\sigma'}(z) \pi_{-p}^\sigma(z) u_p(z) + \\
& \quad \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \sum_{I' \in \text{Succ}_p(I, a)} \pi_p^{\sigma'}(I') v_p^\sigma(I') \\
&= \sum_{z \in Z} \pi_p^{\sigma'}(z) \pi_{-p}^\sigma(z) u_p(z) + \sum_{I \in \mathcal{I}_p} \sum_{I' \in \text{Succ}_p(I)} \pi_p^{\sigma'}(I') v_p^\sigma(I') \\
&= \sum_{z \in Z} \pi_p^{\sigma'}(z) \pi_{-p}^\sigma(z) u_p(z) + \sum_{I \in \mathcal{I}_p: I \neq I(\emptyset)} \pi_p^{\sigma'}(I) v_p^\sigma(I). \tag{A.33}
\end{aligned}$$

According to the definition of expected payoff, $\sum_{z \in Z} \pi_p^{\sigma'}(z) \pi_{-p}^\sigma(z) u_p(z) = v_p^{\langle \sigma'_p, \sigma_{-p} \rangle}$. So,

$$\begin{aligned}
& \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) (v_p^\sigma(I, a) - v_p^\sigma(I)) \sigma'_p(I, a) \\
&= v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} + \sum_{I \in \mathcal{I}_p: I \neq I(\emptyset)} \pi_p^{\sigma'}(I) v_p^\sigma(I) - \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) v_p^\sigma(I) \\
&= v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - v_p^\sigma(I(\emptyset)) \\
&= v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - v_p^\sigma. \tag{A.34}
\end{aligned}$$

So the lemma holds. \square

D. Proof for Lemma 3

Proof. The proof is basically the same as the proof for Lemma 2, except that

$$\begin{aligned}
& \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) v_p^{\prime\sigma}(I, a) \sigma'_p(I, a) \\
&= \sum_{I \in \mathcal{I}_p} \sum_{h \in I} \sum_{a \in A(I)} \sum_{z \in Z: z \in \text{Succ}_p(h, a)} \pi_p^{\sigma'}(z) \pi_{-p}^\sigma(z) u_p(z) + \\
& \quad \sum_{I \in \mathcal{I}_p} \sum_{I' \in \text{Succ}_p(I)} \pi_p^{\sigma'}(I') v_p^{\prime\sigma}(I') \\
&= \sum_{z \in Z} \pi_p^{\sigma'}(z) \pi_{-p}^\sigma(z) u_p(z) + \sum_{I \in \mathcal{I}_p: I \neq I(\emptyset)} \pi_p^{\sigma'}(I) v_p^{\prime\sigma}(I). \tag{A.35}
\end{aligned}$$

Therefore,

$$\begin{aligned}
& \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) (v_p^{\prime\sigma}(I, a) - v_p^{\prime\sigma}(I)) \sigma'_p(I, a) \\
&= v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} + \sum_{I \in \mathcal{I}_p: I \neq I(\emptyset)} \pi_p^{\sigma'}(I) v_p^{\prime\sigma}(I) - \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) v_p^{\prime\sigma}(I) \\
&= v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - v_p^{\prime\sigma}(I(\emptyset)) \\
&= v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - v_p^{\prime\sigma}. \tag{A.36}
\end{aligned}$$

E. Proof for Equation (14)

Proof. Firstly, we have

$$\begin{aligned}
R_p^{T+T'} &= \max_{\sigma'_p \in \Sigma_p} \sum_{t=1}^{T+T'} v_p^{\langle \sigma'_p, \sigma_{-p}^t \rangle} - \sum_{t=1}^{T+T'} v_p^{\sigma^t} \\
&= \left(T v_p^{\prime\sigma} - \sum_{t=1}^T v_p^{\sigma^t} \right) + \\
& \quad \max_{\sigma'_p \in \Sigma_p} \left(\sum_{t=1}^{T+T'} v_p^{\langle \sigma'_p, \sigma_{-p}^t \rangle} - T v_p^{\prime\sigma} - \sum_{t'=1}^{T'} v_p^{\sigma^{t'}} \right). \tag{A.37}
\end{aligned}$$

Recall that we assume $\sigma = \bar{\sigma}^T$. So, for the second term on the right side in the above equation, according to Lemma 2 and Lemma 3, we have

$$\begin{aligned}
& \sum_{t=1}^{T+T'} v_p^{\langle \sigma'_p, \sigma_{-p}^t \rangle} - T v_p^{\prime\sigma} - \sum_{t'=1}^{T'} v_p^{\sigma^{t'}} \\
&= \left(T v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - T v_p^{\prime\sigma} \right) + \sum_{t'=1}^{T'} \left(v_p^{\langle \sigma'_p, \sigma_{-p}^{t'} \rangle} - v_p^{\sigma^{t'}} \right) \\
&= \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) T (v_p^{\prime\sigma}(I, a) - v_p^{\prime\sigma}(I)) \sigma'_p(I, a) + \\
& \quad \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) \sum_{t'=1}^{T'} (v_p^{\sigma^{t'}}(I, a) - v_p^{\sigma^{t'}}(I)) \sigma'_p(I, a) \\
&= \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) \left(R_p^{T'}(I, a) + \sum_{t'=1}^{T'} r_p^{\sigma^{t'}}(I, a) \right) \sigma'_p(I, a). \tag{A.38}
\end{aligned}$$

Note that

$$R_p^{T', T'}(I, a) = R_p^{T'}(I, a) + \sum_{t'=1}^{T'} r_p^{\sigma^{t'}}(I, a). \tag{A.39}$$

So,

$$\begin{aligned}
& \max_{\sigma'_p \in \Sigma_p} \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma'}(I) \left(R_p^{T'}(I, a) + \sum_{t'=1}^{T'} r_p^{\sigma^{t'}}(I, a) \right) \sigma'_p(I, a) \\
&\leq \sum_{I \in \mathcal{I}_p} \max_a (R_p^{T', T'}(I, a))_+, \tag{A.40}
\end{aligned}$$

and

$$R_p^{T+T'} \leq \left(T v_p^{\prime\sigma} - \sum_{t=1}^T v_p^{\sigma^t} \right) + \sum_{I \in \mathcal{I}_p} \max_a (R_p^{T', T'}(I, a))_+. \tag{A.41}$$

\square

F. Proof for Theorem 1

Before we prove the theorem, we would like to introduce a useful lemma.

Lemma A.4. *Given an arbitrary strategy σ , $T > 0$, and $v_p^{\prime\sigma}(I)$ at all infosets, compute $v_p^{\prime\sigma}(I, a)$ according to (8), then, $\max_{\sigma'_p \in \Sigma_p} v_p^{\langle \sigma'_p, \sigma_{-p} \rangle} - \frac{1}{T} \sum_{I \in \mathcal{I}_p} \max_a (R_p^{T'}(I, a))_+ \leq v_p^{\prime\sigma}$.*

\square

Besides, if $v_p^{\prime\sigma}(I) \leq \max_a v_p^{\prime\sigma}(I, a)$ at every infoset, then, $v_p^{\prime\sigma} \leq \max_{\sigma'_p \in \Sigma_p} v_p^{\langle \sigma'_p, \sigma-p \rangle}$.

Proof. For the first inequality, according to Lemma 3,

$$\begin{aligned} & \max_{\sigma'_p \in \Sigma_p} v_p^{\langle \sigma'_p, \sigma-p \rangle} - v_p^{\prime\sigma} \\ &= \max_{\sigma'_p \in \Sigma_p} \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) (v_p^{\prime\sigma}(I, a) - v_p^{\prime\sigma}(I)) \sigma'_p(I, a) \\ &= \max_{\sigma'_p \in \Sigma_p} \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) \frac{1}{T} R_p^{T'}(I, a) \sigma'_p(I, a) \\ &\leq \frac{1}{T} \sum_{I \in \mathcal{I}_p} \max_a (R_p^{T'}(I, a))_+. \end{aligned} \quad (\text{A.42})$$

Rearranging the above equation gives the result.

For the second inequality, let σ_p^* be the strategy that maximizes $\sum_{a \in A(I)} v_p^{\prime\sigma}(I, a) \sigma_p^*(I, a)$ at every infoset, i.e., $\sigma_p^*(I, a) = \mathbb{1}_{a=\arg\max v_p^{\prime\sigma}(I, a)}$. According to Lemma 3,

$$\begin{aligned} & v_p^{\prime\sigma} - v_p^{\langle \sigma_p^*, \sigma-p \rangle} \\ &= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma^*}(I) (v_p^{\prime\sigma}(I) - v_p^{\prime\sigma}(I, a)) \sigma_p^*(I, a) \\ &= \sum_{I \in \mathcal{I}_p} \pi_p^{\sigma^*}(I) \left(v_p^{\prime\sigma}(I) - \max_{a \in A(I)} v_p^{\prime\sigma}(I, a) \right). \end{aligned} \quad (\text{A.43})$$

When $v_p^{\prime\sigma}(I) \leq \max_a v_p^{\prime\sigma}(I, a)$, we have $v_p^{\prime\sigma} \leq v_p^{\langle \sigma_p^*, \sigma-p \rangle} \leq \max_{\sigma'_p \in \Sigma_p} v_p^{\langle \sigma'_p, \sigma-p \rangle}$. \square

Now we can proof Theorem 1.

Proof. According to Lemma A.4, we have

$$\begin{aligned} T v_p^{\prime\sigma} - \sum_{t=1}^T v_p^{\sigma^t} &\leq T \max_{\sigma'_p \in \Sigma_p} v_p^{\langle \sigma'_p, \sigma-p \rangle} - \sum_{t=1}^T v_p^{\sigma^t} \\ &= \max_{\sigma'_p \in \Sigma_p} \sum_{t=1}^T v_p^{\langle \sigma'_p, \sigma-p \rangle} - \sum_{t=1}^T v_p^{\sigma^t} \\ &= R_p^T. \end{aligned} \quad (\text{A.44})$$

Note that $\sigma = \bar{\sigma}^T$. So, according to (14) in the paper,

$$R_p^{T+T'} \leq R_p^T + \sum_{I \in \mathcal{I}_p} \max_a (R_p^{T, T'}(I, a))_+. \quad (\text{A.45})$$

Since $\epsilon(\sigma) = \epsilon(\bar{\sigma}^T) = \frac{1}{T} \sum_{p \in P} R_p^T$ and $\epsilon(\bar{\sigma}^{T, T'}) = \frac{1}{T+T'} \sum_{p \in P} R_p^{T+T'}$, we have

$$\epsilon(\bar{\sigma}^{T, T'}) \leq \frac{T\epsilon(\sigma)}{T+T'} + \frac{\sum_{I \in \mathcal{I}_p} \max_a (R_p^{T, T'}(I, a))_+}{T+T'}. \quad (\text{A.46})$$

According to (12) in the paper, the theorem holds. \square

X. PROOFS FOR RECURSIVE CFR

A. Proof for Proposition 1

We first quote the definition of Generalized Weakened Fictitious Play (GWFP) [47] for completeness. GWFP is a kind of iterative algorithm, as defined in Definition A.2. In

the definition, $\text{BR}_{\epsilon_t}(\Pi_{-p}^t) \in \{\Pi'_p \in \Sigma_p : u_p(\Pi'_p, \Pi_{-p}^t) \geq u_p(\text{BR}(\Pi_{-p}^t), \Pi_{-p}^t) - \epsilon_t\}$ is a ϵ_t -BR against Π_{-p}^t .

Definition A.2. [47] A generalized weakened fictitious play is a process of mixed strategies, $\{\Pi^t\}$, $\Pi^t \in \times_{p \in P} \Sigma_p$, s.t.

$$\Pi_p^{t+1} \in (1 - \alpha_{t+1}) \Pi_p^t + \alpha_{t+1} (\text{BR}_{\epsilon_t}(\Pi_{-p}^t) + M_{t+1}), \forall p \in P,$$

with $\alpha_t \rightarrow 0$ and $\epsilon_t \rightarrow 0$ as $t \rightarrow \infty$, $\sum_{t=1}^{\infty} \alpha_t = \infty$, and $\{M_t\}$ a sequence of perturbations that satisfies $\forall \beta > 0$

$$\lim_{t \rightarrow \infty} \sup_k \left\{ \left\| \sum_{i=t}^{k-1} \alpha_{i+1} M_{i+1} \right\| \text{ s.t. } \sum_{i=t}^{k-1} \alpha_{i+1} \leq \beta \right\} = 0.$$

In [9], a special form of GWFP, named Full-width extensive-form fictitious play (XFP), is given. In XFP, ϵ_t and M_t is set to zero and α_t is set to $\frac{1}{t+1}$ at every iteration. So $\text{BR}_{\epsilon_t}(\Pi_{-p}^t)$ is a best response against the mixed strategy Π_{-p}^t of the opponent, and the mixed strategy Π^t is the average of the best responses.

Proof. ReCFR guarantees that $v_p^{\prime\sigma^t}(I) \leq \max_a v_p^{\prime\sigma^t}(I, a)$ for all infosets. According to Lemma A.4, when $\lambda_p^t(I) = 0$ at every infoset, we have $v_p^{\sigma^{t+1}} = v_p^{\prime\sigma^t} = \max_{\sigma'_p} v_p^{\langle \sigma'_p, \bar{\sigma}^t-p \rangle}$, i.e., σ_p^{t+1} is a best response to $\bar{\sigma}^t-p$. According to the definition of XFP, the proposition holds. \square

B. Proof for Proposition 2

Proof. According to the definition of ReCFR, when $\lambda_p^t(I) = \sum_a (R^t(I, a))_+^2$ and

$$t v_p^{\prime\bar{\sigma}^t}(I, a) = \sum_{k=1}^t v_p^{\sigma^k}(I, a), \quad (\text{A.47})$$

we have

$$t v_p^{\prime\bar{\sigma}^t}(I) = \sum_{k=1}^t v_p^{\sigma^k}(I). \quad (\text{A.48})$$

Note that the two equations also hold when $\sum_a (R^t(I, a))_+^2 = 0$. Then, according to the recursive definition of counterfactual values and RSVs, the above two equations hold at every infoset if $\lambda_p^t(I) = \sum_a (R^t(I, a))_+^2$ at every infoset. Therefore, the substitute regrets recover the cumulative regrets and ReCFR recovers CFR. \square

C. Proof for Theorem 2

Proof. Firstly, for any $\sigma'_p \in \Sigma_p$, we have

$$\begin{aligned} & t v_p^{\prime\bar{\sigma}^t} - (t-1) v_p^{\prime\bar{\sigma}^{t-1}} - v_p^{\sigma^t} \\ &= (t-1) \left(v_p^{\langle \sigma'_p, \bar{\sigma}^{t-1}-p \rangle} - v_p^{\prime\bar{\sigma}^{t-1}} \right) + \\ & \quad \left(v_p^{\langle \sigma'_p, \bar{\sigma}^t-p \rangle} - v_p^{\sigma^t} \right) - t \left(v_p^{\langle \sigma'_p, \bar{\sigma}^t-p \rangle} - v_p^{\prime\bar{\sigma}^t} \right). \end{aligned} \quad (\text{A.49})$$

According to Lemma 3, at iteration $t > 1$, we have

$$\begin{aligned} & t \left(v_p^{\langle \sigma'_p, \bar{\sigma}^t-p \rangle} - v_p^{\prime\bar{\sigma}^t} \right) \\ &= t \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) \left(v_p^{\prime\bar{\sigma}^t}(I, a) - v_p^{\prime\bar{\sigma}^t}(I) \right) \sigma'_p(I, a) \\ &= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) R_p^{t'}(I, a) \sigma'_p(I, a), \end{aligned} \quad (\text{A.50})$$

and

$$\begin{aligned}
& (t-1) \left(v_p^{\langle \sigma'_p, \bar{\sigma}^{t-1} \rangle} - v_p^{\bar{\sigma}^{t-1}} \right) \\
&= (t-1) \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) \\
& \quad \left(v_p^{\bar{\sigma}^{t-1}}(I, a) - v_p^{\bar{\sigma}^{t-1}}(I) \right) \sigma'_p(I, a) \\
&= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) R_p^{t-1}(I, a) \sigma'_p(I, a).
\end{aligned} \tag{A.51}$$

Besides, according to Lemma 2,

$$\begin{aligned}
& v_p^{\langle \sigma'_p, \sigma^t \rangle} - v_p^{\sigma^t} \\
&= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) r_p^{\sigma^t}(I, a) \sigma'_p(I, a).
\end{aligned} \tag{A.52}$$

So,

$$\begin{aligned}
& t v_p^{\bar{\sigma}^t} - (t-1) v_p^{\bar{\sigma}^{t-1}} - v_p^{\sigma^t} \\
&= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) g_p^t(I, a) \sigma'_p(I, a),
\end{aligned} \tag{A.53}$$

where

$$g_p^t(I, a) = R_p^{t-1}(I, a) + r_p^{\sigma^t}(I, a) - R_p^t(I, a). \tag{A.54}$$

Similarly, when $t = 1$, we have

$$\begin{aligned}
& v_p^{\bar{\sigma}^t} - v_p^{\sigma^t} \\
&= (v_p^{\bar{\sigma}^t} - v_p^{\langle \sigma'_p, \bar{\sigma}^t \rangle}) + (v_p^{\langle \sigma'_p, \bar{\sigma}^t \rangle} - v_p^{\sigma^t}) \\
&= \sum_{I \in \mathcal{I}_p} \sum_{a \in A(I)} \pi_p^{\sigma'}(I) \left(r_p^{\bar{\sigma}^t}(I, a) - R_p^t(I, a) \right) \sigma'_p(I, a).
\end{aligned} \tag{A.55}$$

Let $v_p^{\bar{\sigma}^0} = 0$ and $v_p^{\bar{\sigma}^0}(I, a) = v_p^{\bar{\sigma}^0}(I) = R_p^0(I, a) = 0$. Then, (A.53) holds for $t \geq 1$.

Let $\sigma'_p = \sigma_p^{t+1}$. Notice that $\sigma_p^{t+1}(I, a) = \frac{(R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+}$,

$$\begin{aligned}
& \sum_{a \in A(I)} g_p^t(I, a) \sigma_p^{t+1}(I, a) \\
&= \frac{\sum_a (R_p^{t-1}(I, a) - R_p^t(I, a)) (R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+} + \\
& \quad \frac{\sum_a r_p^{\sigma^t}(I, a) (R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+} \\
&\leq \frac{\sum_a ((R_p^{t-1}(I, a))_+ - (R_p^t(I, a))_+) (R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+} + \\
& \quad \frac{\sum_a r_p^{\sigma^t}(I, a) (R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+}.
\end{aligned} \tag{A.56}$$

Besides, we have $\sum_a r_p^{\sigma^t}(I, a) (R_p^{t-1}(I, a))_+ = 0$ when $t > 1$, as $r_p^{\sigma^t}(I, a) = v_p^{\sigma^t}(I, a) - \sum_a v_p^{\sigma^t}(I, a) \sigma^t(I, a)$ and

$\sigma_p^t(I, a) \propto (R_p^{t-1}(I, a))_+$. It is also true when $t = 1$ as $R_p^0(I, a) = 0$. So,

$$\begin{aligned}
& \frac{\sum_a r_p^{\sigma^t}(I, a) (R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+} \\
&= \frac{\sum_a r_p^{\sigma^t}(I, a) ((R_p^t(I, a))_+ - (R_p^{t-1}(I, a))_+)}{\sum_a (R_p^t(I, a))_+} \\
&\leq \frac{\sum_a (r_p^{\sigma^t}(I, a))^2}{2 \sum_a (R_p^t(I, a))_+} + \\
& \quad \frac{\sum_a ((R_p^t(I, a))_+ - (R_p^{t-1}(I, a))_+)^2}{2 \sum_a (R_p^t(I, a))_+}.
\end{aligned} \tag{A.57}$$

The last inequality is derived according to Fenchel-Young inequality. Combine the above two equations, we get

$$\begin{aligned}
& \sum_{a \in A(I)} g_p^t(I, a) \sigma_p^{t+1}(I, a) \\
&\leq \frac{\sum_a ((R_p^{t-1}(I, a))_+ - (R_p^t(I, a))_+) (R_p^t(I, a))_+}{\sum_a (R_p^t(I, a))_+} + \\
& \quad \frac{\sum_a ((R_p^t(I, a))_+ - (R_p^{t-1}(I, a))_+)^2}{2 \sum_a (R_p^t(I, a))_+} + \\
& \quad \frac{\sum_a (r_p^{\sigma^t}(I, a))^2}{2 \sum_a (R_p^t(I, a))_+} \\
&= \frac{\sum_a (R_p^{t-1}(I, a))_+^2 - \sum_a (R_p^t(I, a))_+^2}{2 \sum_a (R_p^t(I, a))_+} + \\
& \quad \frac{\sum_a (r_p^{\sigma^t}(I, a))^2}{2 \sum_a (R_p^t(I, a))_+}.
\end{aligned} \tag{A.58}$$

Notice that $\sum_a (R_p^t(I, a))_+^2 = \lambda_p^t(I)$ ((15) in the paper) and $\sqrt{\sum_a (R_p^t(I, a))_+^2} \leq \sum_a (R_p^t(I, a))_+$, we have

$$\begin{aligned}
& \frac{\sum_a (R_p^{t-1}(I, a))_+^2 - \sum_a (R_p^t(I, a))_+^2}{2 \sum_a (R_p^t(I, a))_+} + \\
&\leq \frac{\lambda_p^{t-1}(I) - \lambda_p^t(I)}{2 \sqrt{\lambda_p^t(I)}}.
\end{aligned} \tag{A.59}$$

Then, according to (A.53),

$$\begin{aligned}
& T v_p^{\bar{\sigma}^T} - \sum_{t=1}^T v_p^{\sigma^t} \\
&= \sum_{t=1}^T \left(t v_p^{\bar{\sigma}^t} - (t-1) v_p^{\bar{\sigma}^{t-1}} - v_p^{\sigma^t} \right) \\
&\leq \sum_{t=1}^T \sum_{I \in \mathcal{I}_p} \frac{\left(\lambda_p^{t-1}(I) - \lambda_p^t(I) + \sum_a (r_p^{\sigma^t}(I, a))^2 \right)_+}{2 \sqrt{\lambda_p^t(I)}}.
\end{aligned} \tag{A.60}$$

Finally, according to (18) in the paper,

$$\begin{aligned}
R_p^T &\leq \left(T v_p^{\bar{\sigma}^T} - \sum_{t=1}^T v_p^{\sigma^t} \right) + \sum_{I \in \mathcal{I}_p} \max_a (R_p^T(I, a))_+ \\
&\leq \sum_{t=1}^T \sum_{I \in \mathcal{I}_p} \frac{\left(\lambda_p^{t-1}(I) - \lambda_p^t(I) + \sum_a (r_p^{\sigma^t}(I, a))^2 \right)_+}{2\sqrt{\lambda_p^t(I)}} + \\
&\quad \sum_{I \in \mathcal{I}_p} \max_a (R_p^T(I, a))_+ \\
&\leq \sum_{t=1}^T \sum_{I \in \mathcal{I}_p} \frac{\left(\lambda_p^{t-1}(I) - \lambda_p^t(I) + \sum_a (r_p^{\sigma^t}(I, a))^2 \right)_+}{2\sqrt{\lambda_p^t(I)}} + \\
&\quad \sum_{I \in \mathcal{I}_p} \sqrt{\lambda_p^T(I)}.
\end{aligned} \tag{A.61}$$

D. Proof for Corollary 1

Proof. Note that

$$\begin{aligned}
\sum_a (r_p^{\sigma^t}(I, a))^2 &\leq (\pi_{-p}^{\sigma^t}(I))^2 \Delta^2(I) |A(I)| \\
&\leq \pi_{-p}^{\sigma^t}(I) \Delta^2(I) |A(I)|.
\end{aligned} \tag{A.62}$$

Note that $\pi_{-p}^{\bar{\sigma}^t}(I)t = \sum_{k=1}^t \pi_{-p}^{\sigma^k}(I)$. When $\lambda_p^t(I) = \pi_{-p}^{\bar{\sigma}^t}(I) \Delta^2(I) |A(I)| t$, we have $\lambda_p^{t-1}(I) - \lambda_p^t(I) + \sum_a (r_p^{\sigma^t}(I, a))^2 \leq 0$. According to Theorem 2,

$$\begin{aligned}
R_p^T &\leq \sum_{I \in \mathcal{I}_p} \sqrt{\lambda_p^T(I)} \\
&\leq \sum_{I \in \mathcal{I}_p} \sqrt{\pi_{-p}^{\bar{\sigma}^T}(I) \Delta(I) \sqrt{|A(I)|} \sqrt{T}}.
\end{aligned} \tag{A.63}$$

As $\epsilon(\bar{\sigma}^T) = \frac{1}{T} \sum_{p \in P} R_p^T$, the corollary holds. \square

E. Proof for Corollary 2

Lemma A.5. [48] *Let $a_t \geq 0$ for $0 \leq t < T$ and $f : [0, +\infty) \rightarrow [0, +\infty)$ a nonincreasing function. Then*

$$\sum_{t=1}^T a_t f \left(a_0 + \sum_{k=1}^t a_k \right) \leq \int_{a_0}^{\sum_{t=0}^T a_t} f(x) dx. \tag{A.64}$$

Proof. [48]. Denote $s_t = \sum_{k=0}^t a_k$.

$$a_t f \left(a_0 + \sum_{k=1}^t a_k \right) = a_t f(s_t) \leq \int_{s_{t-1}}^{s_t} f(x) dx. \tag{A.65}$$

Summing over $t = 1, \dots, T$, we have the stated bound. \square

According to Lemma A.5, if $\sum_{k=1}^t a_k^2 > 0$ for $1 \leq t \leq T$, it is immediately that

$$\sum_{t=1}^T \frac{a_t^2}{\sqrt{\sum_{k=1}^t a_k^2}} \leq 2 \sqrt{\sum_{t=1}^T a_t^2}. \tag{A.66}$$

Now, we can prove the corollary.

Proof. Note that $\pi_{-p}^{\bar{\sigma}^t}(I)t = \sum_{k=1}^t \pi_{-p}^{\sigma^k}(I)$. When $\lambda_p^t(I) = \lambda \pi_{-p}^{\bar{\sigma}^t}(I) \Delta^2(I) |A(I)| t$, we have $\lambda_p^{t-1}(I) \leq \lambda_p^t(I)$. According to Theorem 2,

$$\begin{aligned}
R_p^T &\leq \sum_{I \in \mathcal{I}_p} \left\{ \sum_{t=1}^T \frac{\pi_{-p}^{\sigma^t}(I) \Delta^2(I) |A(I)|}{2\sqrt{\lambda_p^t(I)}} + \sqrt{\lambda_p^T(I)} \right\} \\
&= \sum_{I \in \mathcal{I}_p} \left\{ \sum_{t=1}^T \frac{\pi_{-p}^{\sigma^t}(I) \Delta^2(I) |A(I)|}{2\sqrt{\lambda \sum_{k=1}^t \pi_{-p}^{\sigma^k}(I) \Delta^2(I) |A(I)|}} + \right. \\
&\quad \left. \sqrt{\lambda \pi_{-p}^{\bar{\sigma}^T}(I) \Delta^2(I) |A(I)| T} \right\} \\
&\leq \left(\frac{1}{\sqrt{\lambda}} + \sqrt{\lambda} \right) \sum_{I \in \mathcal{I}_p} \sqrt{\pi_{-p}^{\bar{\sigma}^T}(I) \Delta(I) \sqrt{|A(I)|} \sqrt{T}}.
\end{aligned} \tag{A.67}$$

As $\epsilon(\bar{\sigma}^T) = \frac{1}{T} \sum_{p \in P} R_p^T$, the corollary holds. \square

XI. PROOFS FOR RECURSIVE CFR WITH BOOTSTRAPPING

A. Proof for Equation (23)

Proof. The equation shows that, for any $I \in \mathcal{I}_p, a \in A(I)$,

$$\begin{aligned}
u_p^{\bar{\sigma}^t}(I, a) &= \mathbb{E}_{h \sim I, h' \sim \text{Succ}_p(h \cdot a)} \left\{ \mathbb{1}_{h' \in Z u_p(h')} + \right. \\
&\quad \left. \mathbb{1}_{h' \notin Z u_p^{\bar{\sigma}^t}(I(h'))} \right\}.
\end{aligned} \tag{A.68}$$

The right side of (A.68) can be expanded as

$$\begin{aligned}
&\mathbb{E}_{h \sim I, h' \sim \text{Succ}_p(h \cdot a)} \left\{ \mathbb{1}_{h' \in Z u_p(h')} + \mathbb{1}_{h' \notin Z u_p^{\bar{\sigma}^t}(I(h'))} \right\} \\
&= \sum_{h \in I} p(h|I) \sum_{h' \in \text{Succ}_p(h \cdot a)} p(h'|h \cdot a) \mathbb{1}_{h' \in Z u_p(h')} + \\
&\quad \sum_{h \in I} p(h|I) \sum_{h' \in \text{Succ}_p(h \cdot a)} p(h'|h \cdot a) \mathbb{1}_{h' \notin Z u_p^{\bar{\sigma}^t}(I(h'))},
\end{aligned} \tag{A.69}$$

where $p(h|I)$ is the probability of sampling h when I is reached and $p(h'|h \cdot a)$ is the probability of sampling h' when $h \cdot a$ is reached. Assume player p is using strategy $\hat{\sigma}$, and player $-p$ is using the average strategy $\bar{\sigma}^t$. We have

$$\begin{aligned}
p(h|I) &= \frac{\pi_p^{\hat{\sigma}}(h) \pi_{-p}^{\bar{\sigma}^t}(h)}{\sum_{h \in I} \pi_p^{\hat{\sigma}}(h) \pi_{-p}^{\bar{\sigma}^t}(h)} \\
&= \frac{\pi_{-p}^{\bar{\sigma}^t}(h)}{\sum_{h \in I} \pi_{-p}^{\bar{\sigma}^t}(h)} \\
&= \frac{\pi_{-p}^{\bar{\sigma}^t}(h)}{\pi_{-p}^{\bar{\sigma}^t}(I)}.
\end{aligned} \tag{A.70}$$

Note that $\pi_p^{\hat{\sigma}}(h) = \pi_p^{\hat{\sigma}}(I)$ for any $h \in I$. As for $p(h'|h \cdot a)$, we have

$$\begin{aligned}
p(h'|h \cdot a) &= \frac{\pi_p^{\hat{\sigma}}(h') \pi_{-p}^{\bar{\sigma}^t}(h')}{\pi_p^{\hat{\sigma}}(h \cdot a) \pi_{-p}^{\bar{\sigma}^t}(h \cdot a)} \\
&= \pi_{-p}^{\bar{\sigma}^t}(h \cdot a, h')
\end{aligned} \tag{A.71}$$

Note that $\pi_p^{\hat{\sigma}}(h') = \pi_p^{\hat{\sigma}}(h \cdot a)$ as h' is the earliest reachable history of player p from $h \cdot a$. Put them into (A.69), we get

$$\begin{aligned}
& \mathbb{E}_{h \sim I, h' \sim \text{Succ}_p(h \cdot a)} \left\{ \mathbb{1}_{h' \in Z} u_p(h') + \mathbb{1}_{h' \notin Z} u_p^{\bar{\sigma}^t}(I(h')) \right\} \\
&= \frac{1}{\pi_{-p}^{\bar{\sigma}^t}(I)} \sum_{h \in I} \sum_{h' \in Z: h' \in \text{Succ}_p(h \cdot a)} \pi_{-p}^{\bar{\sigma}^t}(h) \pi_{-p}^{\bar{\sigma}^t}(h \cdot a, h') u_p(h') + \\
& \quad \frac{1}{\pi_{-p}^{\bar{\sigma}^t}(I)} \sum_{h \in I} \sum_{h' \notin Z: h' \in \text{Succ}_p(h \cdot a)} \pi_{-p}^{\bar{\sigma}^t}(h) \pi_{-p}^{\bar{\sigma}^t}(h \cdot a, h') u_p^{\bar{\sigma}^t}(I(h')) \\
&= \frac{1}{\pi_{-p}^{\bar{\sigma}^t}(I)} \sum_{h \in I} \sum_{h' \in Z: h' \in \text{Succ}_p(h \cdot a)} \pi_{-p}^{\bar{\sigma}^t}(h') u_p(h') + \\
& \quad \frac{1}{\pi_{-p}^{\bar{\sigma}^t}(I)} \sum_{I' \in \text{Succ}_p(I, a)} v_p^{\bar{\sigma}^t}(I'), \\
&= \frac{1}{\pi_{-p}^{\bar{\sigma}^t}(I)} v_p^{\bar{\sigma}^t}(I, a). \tag{A.72}
\end{aligned}$$

The second equality is derived according to the the definition of $\text{Succ}_p(I, a)$. The last equality is because of the recursive property of RSVs. Because $u_p^{\bar{\sigma}^t}(I, a) = v_p^{\bar{\sigma}^t}(I, a) / \pi_{-p}^{\bar{\sigma}^t}(I)$, the equation holds. \square

B. Proof for Theorem 3

The Bootstrap learning at every iteration in the algorithm is a mimic of Q-learning. According to [49], the convergence of Q-learning is guaranteed by Theorem A.4.

Theorem A.4. [49] *The Q-learning algorithm given by*

$$\begin{aligned}
Q_{k+1}(s, a) &= Q_k(s, a) + \\
& \alpha_k \left(r(s, a) + \gamma \max_{a'} Q_k(s', a') - Q_k(s, a) \right), \tag{A.73}
\end{aligned}$$

where s' is the state transferred from s after selecting action a , converges to the optimal $Q^*(s, a)$ values if

- The state and action spaces are finite.
- $\sum_k \alpha_k = \infty$ and $\sum_k \alpha_k^2 < \infty$ uniformly w.p.1.
- $\text{Var}(r(s, a))$ is bounded.

Similarly, we can prove Theorem 3 in the paper.

Lemma A.6. *For any $x \in \mathbb{R}^n$, $y \in \mathbb{R}^n$, $X \in \mathbb{R}$ and $Y \in \mathbb{R}$, if $\sum_{i=1}^n (x_i - X)_+^2 = \sum_{i=1}^n (y_i - Y)_+^2 > 0$, then, $|X - Y| \leq \max_i |x_i - y_i|$.*

Proof. Firstly, for any $z \in \mathbb{R}$, we have

$$\sum_{i=1}^n (x_i - X)_+^2 = \sum_{i=1}^n ((y_i + z) - (Y + z))_+^2. \tag{A.74}$$

Let $z = \max_i |x_i - y_i|$. Note that $x_{i'} \leq y_{i'} + \max_i |x_i - y_i|$ for any $1 \leq i' \leq n$. So, we have $X \leq Y + \max_i |x_i - y_i|$. Otherwise, (A.74) is contradicted. Similarly, we have $Y \leq X + \max_i |x_i - y_i|$. So, the lemma holds. \square

Now, we are ready to prove Theorem 3 in the paper.

Proof. At infoset I at iteration t , let a state s in Theorem A.4 represent an infoset I in Recursive CFR. Let $Q_{k+1}(s, a) = u_p^{\bar{\sigma}^t, k+1}(I, a)$, $Q_k(s', a') = u_p^{\bar{\sigma}^t, k}(I', a')$, and $r(s, a) = \mathbb{E}_{h \sim I, h' \sim \text{Succ}_p(h \cdot a)} \mathbb{1}_{h' \in Z} u_p(h')$. Following the proof in [49],

we only need to prove that the following mapping operator is a contraction operator:

$$\mathbf{H} : (\mathbf{H}u_p^{\bar{\sigma}^t, k})(I, a) = r(I, a) + \gamma \sum_{I' \in \mathcal{I}_p} P_p^{\bar{\sigma}^t}(I'|I, a) u_p^{\bar{\sigma}^t, k}(I') \tag{A.75}$$

where $P_p^{\bar{\sigma}^t}(I'|I, a) = \mathbb{1}_{I' \in \text{Succ}_p(I, a)} \pi_{-p}^{\bar{\sigma}^t}(I') / \pi_{-p}^{\bar{\sigma}^t}(I)$ is the probability of reach infoset I' from I . So we need to prove

$$\left\| (\mathbf{H}u_p^{\bar{\sigma}^t, k}) - (\mathbf{H}u_p^{\bar{\sigma}^t, k'}) \right\|_{\infty} \leq \gamma \left\| u_p^{\bar{\sigma}^t, k} - u_p^{\bar{\sigma}^t, k'} \right\|_{\infty}, \tag{A.76}$$

for any $\bar{\sigma}^t$ and $1 \leq k, k' \leq K$. Since $r(I, a)$ is constant with respect to k ,

$$\begin{aligned}
& \left\| (\mathbf{H}u_p^{\bar{\sigma}^t, k}) - (\mathbf{H}u_p^{\bar{\sigma}^t, k'}) \right\|_{\infty} \\
&= \gamma \max_{I, a} \left| \sum_{I' \in \mathcal{I}_p} P_p^{\bar{\sigma}^t}(I'|I, a) \left(u_p^{\bar{\sigma}^t, k}(I') - u_p^{\bar{\sigma}^t, k'}(I') \right) \right|. \tag{A.77}
\end{aligned}$$

Since $P_p^{\bar{\sigma}^t}(I'|I, a) \geq 0$ and $\sum_{I' \in \mathcal{I}_p} P_p^{\bar{\sigma}^t}(I'|I, a) \leq 1$,

$$\begin{aligned}
& \left\| (\mathbf{H}u_p^{\bar{\sigma}^t, k}) - (\mathbf{H}u_p^{\bar{\sigma}^t, k'}) \right\|_{\infty} \\
&\leq \gamma \max_I \left| u_p^{\bar{\sigma}^t, k}(I) - u_p^{\bar{\sigma}^t, k'}(I) \right| \\
&\leq \gamma \max_I \max_a \left| u_p^{\bar{\sigma}^t, k}(I, a) - u_p^{\bar{\sigma}^t, k'}(I, a) \right| \\
&= \gamma \left\| u_p^{\bar{\sigma}^t, k} - u_p^{\bar{\sigma}^t, k'} \right\|_{\infty}. \tag{A.78}
\end{aligned}$$

The second inequality is because of Lemma A.6. Note that when $\lambda_p^t(I) = 0$, i.e., $u_p^{\bar{\sigma}^t, k}(I) = \max_a u_p^{\bar{\sigma}^t, k}(I, a)$ and $u_p^{\bar{\sigma}^t, k'}(I) = \max_a u_p^{\bar{\sigma}^t, k'}(I, a)$, the inequality also holds. So, when $\gamma < 1$, \mathbf{H} is a γ -contraction mapping. According to the updating rules of RSVs ((22) and (23) in the paper), $u_p^{\bar{\sigma}^t}(I, a)$ is a fixed point of the mapping shown in (A.75). So, $u_p^{\bar{\sigma}^t, K}(I, a)$ converges to $u_p^{\bar{\sigma}^t}(I, a)$ w.p.1 when $K \rightarrow \infty$ for any (I, a) . Moreover, when $\gamma = 1$, according to [49], $u_p^{\bar{\sigma}^t, K}(I, a)$ can still converge, as long as all the terminal histories are visited with probabilities greater than 0. Since we only consider depth-limited games, the theorem holds. \square

XII. RULES OF HEADS-UP LIMIT TEXAS HOLD'EM AND FLOP HOLD'EM POKER

Heads-up Limit Texas Hold'em (HULH) is a two-player zero-sum game. At the beginning of the game, player 1 should place \$50 and player 2 should place \$100 on the desk. Then 2 private cards are dealt to each player and the game goes to the first round. The game contains 4 rounds. In each round, the two players take action in turn. A player can choose fold, call or raise. Action fold means the player gives up the game thus the game terminates immediately and the player loses the money on the desk. If a player chooses to call, he should place the same money to the desk as the other player, and the game goes to the next round or terminates if it is the final round. When a player chooses to raise, he should place more money on the desk than the other player. However, there can not be more than three raises in the first two rounds and more than four raises in the second two rounds. Raises in the first two rounds are \$100 and raises in the second two rounds are

\$200. In the first round, player 1 should act first, while player 2 acts first in the rest of the rounds. When the first round ends, three community cards are dealt face-up on the desk and the second round starts. Another two community cards are dealt at the beginning of the next two rounds, one for each round, followed by a series of betting. If the rounds end without any player folds, the two players should reveal their private cards and the five cards for the players (2 private cards plus three community cards) are compared. Then the money on the desk is won by the player with stronger cards or split evenly if a tie.

Flop Hold'em Poker (FHP) is a simplified HULH, which only contains the first two rounds of betting.

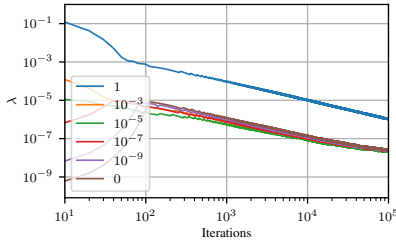


Fig. 9: Curves of the adaptive λ in ReCFR with different initial values on Leduc poker.

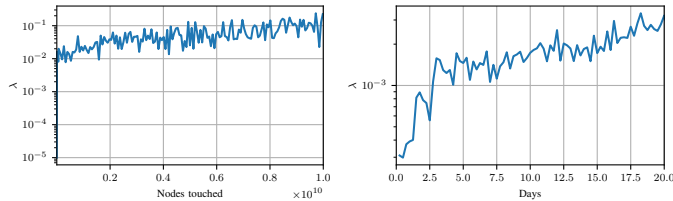


Fig. 10: Curves of the adaptive λ in Neural ReCFR-B on FHP (left) and HULH poker (right).

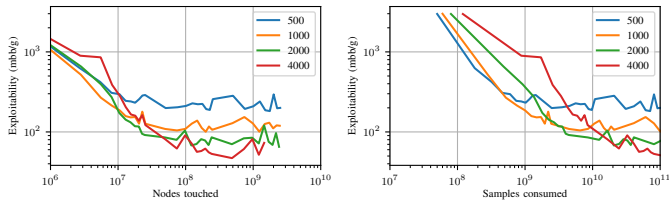


Fig. 11: Exploitability curves of Deep CFR with different SGD steps on FHP. The x-axes represent the number of nodes touched and the number of samples consumed, respectively.

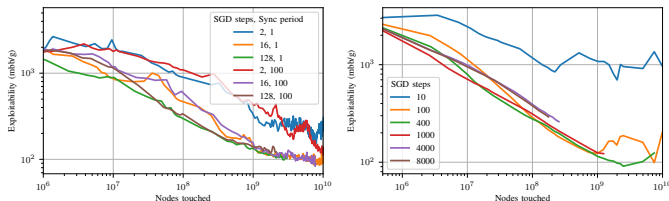


Fig. 12: Exploitability curves of NFSP and DNCFR with different SGD steps.

XIII. HYPER-PARAMETERS AND EXPERIMENTAL ENVIRONMENT

The hyper-parameters for all the algorithms are listed in Table I. We implement Deep CFR and DREAM with the default hyper-parameters given in [18] and [20], respectively. The hyper-parameters for NFSP and DNCFR are decided according to a set of experiments, see Figure 12.

All the algorithms are implemented in C++ based on TensorFlow [50] C++ API. All the experiments are run in a high performance computing cluster using 10 CPU cores and 40GB of memory for 10 days. The experiments for Neural CFR-B and Deep CFR on HULH are conducted on a server with 10 CPU cores, 100GB of memory and one 2080TI GPU. The operating system is Ubuntu 18.04 and the compiler is GCC-9.0. The implementations are based on an open-source framework “OpenSpiel”[39]. We seed the random generator using a real random device (std::random_device in C++).

XIV. ADDITIONAL RESULTS

In this section, additional results of ReCFR, Neural ReCFR-B, and other algorithms are given.

In Figure 9 and 10, the curves of the adaptive λ in ReCFR and Neural ReCFR-B are given. As we can see, the λ will generally converge, and it converges to different values for different games.

In Figure 11, the results of Deep CFR with different SGD steps are given. As we can see, Deep CFRs with fewer SGD steps converge earlier to higher exploitability. The results of DNCFR and NFSP on FHP are given in Figure 12.

TABLE I: Hyper-parameters

Algorithm	Hyper-parameters
Deep CFR	sampling method = external sampling CFR [33], traversals per iteration = 10,000, regret memory size = 40 million, strategy memory size = 40 million, optimizer = Adam, learning rate = 0.001, batch size = 10,000, SGD steps per training for regret network = 4,000, SGD steps per training for strategy network = 4,000, training regret network from scratch = true, training strategy network from scratch = false.
Deep OSCFR	Same as Deep CFR except sampling method = outcome sampling CFR [33], traversals per iteration = 50,000.
Deep CFR on HULH	Same as Deep CFR except batch size = 20,000, SGD steps per training for regret network = 32,000, SGD steps per training for strategy network = 32,000.
DREAM	Same as Deep CFR except sampling method = outcome sampling CFR [33], traversals per iteration = 50,000, SGD steps per training for regret network = 3,000, SGD steps per training for strategy network = 3,000, global value memory size = 200,000, batch size for global value network = 512, SGD steps per training for global value network = 1,000.
DNCFR	sampling method = CFR+ with robust sampling [19], traversals per iteration = 10,000, regret memory = none, strategy memory = none, optimizer = Adam, learning rate = 0.001, batch size = 10,000, SGD steps per training for regret network = 400, SGD steps per training for strategy network = 400, training regret network from scratch = false, training strategy network from scratch = false.
NFSP	sampling method = self-play (trajectory sampling), plays per iteration = 1,000, value memory size = 1 million, strategy memory size = 10 million, optimizer = Adam, learning rate = 0.001, batch size = 128, SGD steps per training for value network = 16, SGD steps per training for strategy network = 16, training value network from scratch = false, training strategy network from scratch = false, anticipatory parameter = 0.1.
Neural CFR-B on HULH	Same as NFSP except RSV memory size = none (default) / 1 million (with memory), SGD steps per training for RSV network = 2 epochs (default) / 32 (with memory), training RSV network from scratch = false, $\lambda_{init} = 10^{-5}$, $\beta_{amp} = 1.01$, $\beta_{damp} = 0.99$.
Neural CFR-B on HULH	Same as Neural CFR-B except plays per iteration = 100,000, RSV memory size = 4 million, strategy memory size = 40 million, batch size = 6400, SGD steps per training for RSV network = 64, SGD steps per training for strategy network = 64, $\lambda_{init} = 10^{-4}$.