# Pure Monte Carlo Counterfactual Regret Minimization

**Qi Ju** [1], **Ting Feng** [1], **Falin Hei** [1], **Zhemei Fang** [1], **Yunfeng Luo** [1]

[1]Huazhong University of Science and Technology

juqi@hust.edu.cn, fengting@hust.edu.cn, heifalin@hust.edu.cn, zmfang2018@hust.edu.cn, yfluo@hust.edu.cn

## Abstract

Counterfactual Regret Minimization (CFR) and its variants are the best algorithms so far for solving large-scale incomplete information games. Building upon CFR, this paper proposes a new algorithm named Pure CFR (PCFR) for achieving better performance. PCFR can be seen as a combination of CFR and Fictitious Play (FP), inheriting the concept of counterfactual regret (value) from CFR, and using the best response strategy instead of the regret matching strategy for the next iteration. Our theoretical proof that PCFR can achieve Blackwell approachability enables PCFR's ability to combine with any CFR variant including Monte Carlo CFR (MC-CFR). The resultant Pure MCCFR (PMCCFR) can significantly reduce time and space complexity. Particularly, the convergence speed of PMCCFR is at least three times more than that of MCCFR. In addition, since PMCCFR does not pass through the path of strictly dominated strategies, we developed a new warm-start algorithm inspired by the strictly dominated strategies elimination method. Consequently, the PMCCFR with new warm start algorithm can converge by two orders of magnitude faster than the CFR+ algorithm.

## Introduction

A game simulates the strategic interactions between players, and a game problem aims to find the Nash Equilibrium (NE) in which no player can improve by deviating from the equilibrium. While the complete information games have been well addressed by the new artificial intelligence algorithms (Silver et al. 2018), the extensive form games with incomplete information, commonly appearing in the areas of video game design, electricity market, and advertising auction, still face significant challenges. In a complete information game, we can decompose the game into sub-games, and use the backward induction algorithm to solve the equilibrium. However, the hidden information for each player prevents the direct use of sub-game strategy in backward induction algorithm, which largely increases the time and space complexity of solving the equilibrium.

The Counterfactual Regret Minimization (CFR) algorithm (Zinkevich et al. 2007) improved from Regret Matching (RM) (Hart and Mas-Colell 2000) is the basic method for solving incomplete information extensive form games. However, CFR has to traverse the entire game tree in one iteration. With current hardware technology, it can only be used to solve the heads-up limit Texas hold'em poker, and

its information set is $10^{14}$. The more popular no-limit Texas hold'em has $10^{162}$ information sets (Moravík et al. 2017; Johanson 2013), DouDizhu has $10^{83}$ information sets (Zha et al. 2019, 2021), and Mahjong has $10^{121}$ information sets (Li et al. 2020), so it is unrealistic to traverse all information sets in one iteration. In order to overcome the problem that the vanilla CFR rely on game tree traversal, Monte Calo CFR (MCCFR) (Lanctot et al. 2010) provides a sampling-based algorithm to handle the extensive form games. Existing experiments show that the convergence speed of MCCFR is similar to that of CFR. So MCCFR is the preferred algorithm for solving large-scale practical problems. In addition to MCCFR, pruning is also a simple and efficient method to reduce the space complexity of the algorithm. The simplest pruning method is to directly skip the sub-game tree with probability of 0. This pruning alone may avoid 96% of invalid game tree traversal (Lanctot et al. 2010).

In order to increase the convergence rates of CFR, CFR+ (Tammelin 2014) was improved. Although theoretical understanding of their success in practice is still a mystery (Farina et al. 2023), in practice CFR+ converges much faster than vanilla CFR. CFR+ was used to solve heads-up limit Texas hold'em poker (Bowling et al. 2015) and heads-up no-limit Texas hold'em poker (Brown and Sandholm 2019b). Warm start is also a technique to accelerate convergence. The idea is that many games (such as Texas hold'em and Go) have been deeply studied by humans, or have approximated solutions in early small-scale training, while CFR and its variants converge slowly in the early stages of training. If we can start with the approximate solution to the game, we may find equilibrium faster by skipping the early stage of training.

Fictitious Play (FP) is another type of methods for solving NE, which can be traced back to Brown's article in 1951 (Berger 2007; Brown 1951). In *The Theory of Learning in Game* (Fudenberg and Levine 1998) summarized the previous research and defined the form of FP canonically. Generalized weakened fictitious play (GWFP) process (Leslie and Collins 2006) strictly proved that in the presence of certain disturbances and errors, GWFP will eventually converge to NE like FP. Hendon et al. (Hendon, Jacobsen, and Sloth 1996) extended FP to extensive form games. Full-width extensive-form fictitious play (XFP) algorithm (Hein-

rich, Lanctot, and Silver 2015), basis on GWFP, enables FP to converge faster in extensive form games. However, compared with RM, the FP algorithm lags behind in the number of papers and engineering implementation.

Our work can be seen as a combination of MCCFR and FP. From the perspective of MCCFR, our work adopts the Best Response (BR) strategy instead of the RM strategy as the strategy for the next iteration. From the perspective of FP, we introduce the concept of counterfactual regret (value) into calculation of BR strategy, which greatly facilitates the use of FP algorithm in extensive form games. Finally, PMC-CFR converges 3 to 4 times faster than MCCFR in all our experiments, and in some specific problems, the convergence speed can be increased by two orders of magnitude compared with CFR+.

## Notation and Preliminaries

### Game and Blackwell Approachability

**Normal Form Games**  Normal form game is the most basic game Theory model. $\mathcal{N} = \{1, 2, \ldots, i, \ldots\}$ denotes the set of players in a game. Player $i$ has a finite pure strategy set $\mathcal{A}^i$, and $l^i = |\mathcal{A}^i|$ indicates the number of $i$'s pure strategy (where $|\cdot|$ represents the number of elements in the set). A mixed strategy set $\Sigma^i \in \mathbb{R}^{l^i}$ is defined as the set of probability distributions over $\mathcal{A}^i$. A strategy profile $\sigma = \times_{i=1}^{N} \sigma^i$ is a collection of strategies for all players, $\sigma^{-i}$ refers to all strategies in $\sigma$ except $\sigma^{-i}$. We write $u^i(a, \sigma^{-i})$ (respectively $u^i(\sigma^i, \sigma^{-i})$) for the expected reward to Player $i$ if they select pure strategy $a$ (respectively mixed strategy $\sigma^i$) and all other players play the mixed strategy profile $\sigma^{-i}$.

The set of best responses of player $i$ to their oppenents' strategies $\sigma^{-i}$ is

$$b^i(\sigma^{-i}) = \arg\max_{\sigma^{i*} \in \Sigma^i} u^i(\sigma^{i*}, \sigma^{-i}) \tag{1}$$

This $\arg\max$ refers to returning the element with the max value in the set, and returning one of them randomly if there are multiple identical max value. Define the player $i$'s exploitability $\epsilon^i$ for strategy profile $\sigma$ as $\epsilon^i(\sigma) = u^i(b^i(\sigma^{-i}), \sigma^{-i}) - u^i(\sigma)$, and total exploitability as $\epsilon(\sigma) = \sum_{i \in \mathcal{N}} \epsilon^i(\sigma)$. A Nash equilibrium of a game is a strategy profile $\sigma$ that satisfies $\epsilon(\sigma) = 0$.

**Extensive Form Games**  Extensive form game is generally represented as game trees and consist of the following elements:

- $\mathcal{N} = \{1, 2, \ldots i \ldots\}$ represents a collection of players.

- The nodes in the game tree represent possible states in a game $s$, and these nodes constitute a set of states $s \in \mathcal{S}$. The leaf nodes of the game tree $z \in \mathcal{Z}$ are also called terminal states.

- For each state $s \in \mathcal{S}$, its subsequent edges define the state $s$ set of actions that the player or chance in the state can take $A(s)$. $P : \mathcal{S} \to \mathcal{N} \cup \{c\}$ is the player function determines who takes an action in a given state. If $P(s) = c$ then chance determines the action taken in state $s$.

- In a game, players may only know that they are in a certain type of states, but they cannot determine which specific state $s$ they are in. Define the information set $I \in \mathcal{I}^i$ to represent the set of states that the player $i$ cannot distinguish.

- Define a payout function that $R : \mathcal{Z} \to \mathbb{R}^{|\mathcal{N}|}$ maps the end point state to a vector whose components correspond to the payoff for each player.

- The behavioral strategy $\sigma^i(I) \in \mathbb{R}^{|\mathcal{A}(I)|}$ for all $I \in \mathcal{I}^i$ of extensive games is defined independently on each information set.

### Blackwell Approachability Game

**Definition 1**  *A Blackwell approachability game in normal form two-player games can be described as a tuple $(\Sigma, u, S^1, S^2)$, where $\Sigma$ is a strategy profile, $u$ is the payoff function, and $S^i = \mathbb{R}^{l^i}_{\leq 0}$ is a closed convex target cone. The Player $i$'s regret vector of the strategy profile $\sigma$ is $R^i(\sigma) \in \mathbb{R}^{l^i}$, for each component $R^i(\sigma)(a_x) = u^i(a_x, \sigma^{-i}) - u^i(\sigma)$, $a_x \in \mathcal{A}^i$ the average regret vector for players $i$ to take actions at $T$ time $a$ is $\bar{R}^i_T$*

$$\bar{R}^i_T = \frac{1}{T} \sum_{t=1}^{T} R^i(\sigma_t) \tag{2}$$

*At each time $t$, the two players interact in this order:*

- *Player 1 chooses a strategy $\sigma^1_t \in \Sigma^1$;*

- *Player 2 chooses an action $\sigma^2_t \in \Sigma^2$, which can depend adversarially on all the $\sigma^t$ output so far;*

- *Player 1 gets the vector value payoff $R^1(\sigma_t) \in \mathbb{R}^{l^1}$.*

*The goal of Player 1 is to select actions $\sigma^1_1, \sigma^1_2, \ldots \in \Sigma^1$ such that no matter what actions $\sigma^2_1, \sigma^2_2, \ldots \in \Sigma^2$ played by Player 2, the average payoff vector converges to the target set $S^1$.*

$$\min_{\hat{s} \in S^1} \|\hat{s} - R^1_T\|_2 \to 0 \quad as \quad T \to \infty \tag{3}$$

Before explaining how to choose the action $\sigma_t$ to ensure this goal achieve, we first need to define the forceable half-space:

**Definition 2**  *Let $\mathcal{H} \subseteq \mathbb{R}^d$ as half-space, that is, for some $\boldsymbol{a} \in \mathbb{R}^d$, $b \in \mathbb{R}$, $\mathcal{H} = \{\boldsymbol{x} \in \mathbb{R}^d : \boldsymbol{a}^\top \boldsymbol{x} \leq b\}$. In Blackwell approachability games, the halfspace $\mathcal{H}$ is said to be forceable if there exists a strategy $\sigma^{i*} \in \Sigma^i$ of Player $i$ that guarantees that the regret vector $R^i(\sigma)$ is in $\mathcal{H}$ no matter the strategy played by Player $-i$, such that*

$$R^i(\sigma^{i*}, \hat{\sigma}^{-i}) \in \mathcal{H} \quad \forall \hat{\sigma}^{-i} \in \Sigma^{-i} \tag{4}$$

*And $\sigma^{i*}$ is forcing action for $\mathcal{H}$.*

Blackwell's approachability theorem states the following.

**Theorem 1 (Blackwell's theorem)**  *Goal (3) can be attained if and only if every halfspace $\mathcal{H}_t \supseteq S$ is forceable.*

The relationship between Blackwell approachability and no-regret learning is:

**Theorem 2** *Any strategy (algorithm) that achieves Black-well approachability can be converted into an algorithm that achieves no-regret, and vice versa (Abernethy, Bartlett, and Hazan 2011)*

Let $\bar{\sigma}_T^i$ be the average strategy of player $i$:

$$\bar{\sigma}_T^i(a) = \frac{\sum_{t=1}^T \sigma_t^i(a)}{T} \tag{5}$$

In a two-player zero-sum game, the exploitability of the average strategy $\bar{\sigma}_T^i$ at time $T$ of player $i$ is $\epsilon_T^i = \max_{a' \in \mathcal{A}^i} \bar{R}_T^i(a')$ (Brown 2020). Obviously, the regret value is always greater than the exploitability:

$$\lim_{T \to \infty} \epsilon_T = \lim_{T \to \infty} \sum_{i \in \mathcal{N}} \epsilon_T^i \leq \lim_{T \to \infty} \sum_{i \in \mathcal{N}} \min_{\hat{s} \in S^i} \|\hat{s} - R_T^i\|_2 = 0 \tag{6}$$

So, if the algorithm achieves Blackwell approachability, the average strategy $\bar{\sigma}_T^i$ will converge to equilibrium with $T \to \infty$. The rate of convergence is $\epsilon_T^i \leq \bar{R}_T^i \leq L\sqrt{|\mathcal{A}|}/\sqrt{T}$, where $L = \max_{\sigma \in \Sigma, i \in \mathcal{N}} u^i(\sigma) - \min_{\sigma \in \Sigma, i \in \mathcal{N}} u^i(\sigma)$ represents the payoff interval of the game.

## Some basic game solving algorithms

**FP and GWFP** In a vanilla FP process of the normal form game, assuming all players start with random strategy profile $\sigma_{t=1}$, the strategy profile is updated following the function:

$$\sigma_{t+1} = \left(1 - \frac{1}{t+1}\right)\sigma_t + \frac{1}{t+1}b(\sigma_t) \tag{7}$$

Where $t$ represents the number of iterations. When $t \to \infty$, $\sigma_t$ converges to NE.

GWFP is a process of strategy profile $\{\sigma_t\}_{t \geq 0}$, s.t.

$$\sigma_{t+1} = (1 - \alpha_{t+1})\sigma_t + \alpha_{t+1}(b_{\epsilon_t}(\sigma_t) + M_{t+1}) \tag{8}$$

With $\lim_{t \to \infty} \alpha_t = 0$, $\lim_{t \to \infty} \epsilon_t = 0$, $\sum_{t \geq 1} \alpha_t = \infty$, and $\{M_t\}_{t \geq 1}$ a sequence of perturbations that satisfies $\forall E > 0$:

$$\lim_{t \to \infty} \sup_k \left\{ \left\|\sum_{j=t}^{k-1} \alpha_{j+1}M_{j+1}\right\| : \sum_{j=t}^{k-1} \alpha_{j+1}E \right\} = 0 \tag{9}$$

When these three conditions are met, $\sigma_t$ converges to NE with $t \to \infty$ as vanilla FP. Vanilla FP can be seen as a GWFP with stepsize $\alpha_t = \frac{1}{t}$ and $M_t = \mathbf{0}$. $\alpha_t = \frac{1}{\log(t+c)}$, $\alpha_t = \frac{2}{t+1}$ are also common stepsize choices.

**RM and CFR** The RM algorithm satisfies Blackwell approachability that chooses the strategy for the next iteration by:

$$\sigma_{t+1}^i(a) = \begin{cases} \frac{\bar{R}_t^{i,+}(a)}{\sum_{a' \in \mathcal{A}^i} \bar{R}_t^{i,+}(a')} & \text{if } R_t^{i,+} \neq \mathbf{0} \\ 1/l^i & \text{otherwise.} \end{cases} \tag{10}$$

Where $\bar{R}_t^{i,+}(a) = \max(\bar{R}_t^i(a), 0)$. Since the probability of taking action $\sigma_{t+1}^i(a)$ is proportional to the regret value $\bar{R}_t^{i,+}(a)$ of this action, so this algorithm is called the regret matching algorithm, and $\sigma_{t+1}^i$ is called the regret matching strategy. If all players adopt a regret matching strategy $\sigma_{t+1}^i$, then the forceable half space at time $t$ is $\mathcal{H}_t^i := \left\{z \in \mathbb{R}^{l^{-i}} : \langle \bar{R}_t^{i,+}, z \rangle \leq 0\right\}$. Obviously $\mathcal{H}_t^i \supseteq S^i$, so $\bar{\sigma}_t$ can converge to NE with $t \to \infty$.

In an extensive form game, we first consider a particular set of information $I \in \mathcal{I}^i$ and player $i$'s choices made in that information set. Define counterfactual value $u(I, \sigma)$ to be the expected value given that information set $I$ is reached and all players play using strategy $\sigma$ except that player $i$ plays to reach $I$. Finally, for all $a \in A^i(I)$, define $\sigma|_{I \to a}$ to be a strategy profile identical to $\sigma$ except that player $i$ always chooses action $a$ when in information set $I$. The immediate counterfactual regret is:

$$R_{T,\text{imm}}^i(I, a) = \frac{1}{T}\sum_{t=1}^T \pi_{\sigma_t}^{-i}(I) \left(u^i(I, \sigma_t|_{I \to a}) - u^i(I, \sigma_t)\right) \tag{11}$$

$\pi_{\sigma_t}^{-i}(I)$ is the probability of information set $I$ occurring if all players (including chance, except $i$) choose actions according to $\sigma_t$. Let $\bar{R}_{T,\text{imm}}^{i,+}(I, a) = \max(\bar{R}_{T,\text{imm}}^i(I, a), 0)$, the strategy at time $T + 1$ is:

$$\sigma_{T+1}^i(I, a) = \begin{cases} \frac{R_{T,\text{imm}}^{i,+}(I,a)}{\sum_{a \in A(I)} R_{T,\text{imm}}^{i,+}(I,a)} & \text{if } R_{T,\text{imm}}^{i,+}(I) \neq \mathbf{0} \\ \frac{1}{|A(I)|} & \text{otherwise.} \end{cases} \tag{12}$$

The average strategy $\bar{\sigma}_T^i$ for an information set $I$ on iteration $T$ is:

$$\bar{\sigma}_T^i(I) = \frac{\sum_{t=1}^T \pi_{\sigma_t}^i(I)\sigma_t^i(I)}{\sum_{t=1}^T \pi_{\sigma_t}^i(I)} \tag{13}$$

Eventually, $\bar{\sigma}_T$ will converge to NE with $T \to \infty$.

## Variations of CFR

**CFR+ and Weighted Averaging Schemes for CFR** CFR+ is like CFR but with the following small changes. First, after each iteration any action with negative regret is directly set to zero regret

$$\bar{R}_{T,\text{imm}}^{i,+}(I) = \max\left(\frac{T-1}{T}\bar{R}_{T-1,\text{imm}}^{i,+}(I) + \frac{1}{T}R_{T,\text{imm}}^{i,+}(I), 0\right) \tag{14}$$

Second, like different stepsize $\alpha_t$ can be selected in FP, different average weights $w_t$ can be selected for the CFR (they are essentially the same). Let $w_t$ be the weight of strategy $\sigma_t$ s.t. $w_t > 0$ and $w_i \leq w_j \forall i < j$, then the weighted average strategy is defined as

$$\bar{\sigma}_T^{w,i}(I) = \frac{\sum_{t \in T}\left(w_t \pi_{\sigma_t}^i(I)\sigma_i^t(I)\right)}{\sum_{t \in T}\left(w_t \pi_i^{\sigma^t}(I)\right)} \tag{15}$$

will also converge to NE. The convergence rate of this averaging scheme $w$ is:

$$R_T^{i,w} = \max_{a \in A} \frac{\sum_{t=1}^T\left(w_t R_{t,\text{imm}}^{i,w}(a)\right)}{\sum_{t=1}^T w_t} \leq \frac{L\sqrt{|A|}\sqrt{\sum_{t=1}^T w_t^2}}{\sum_{t=1}^T w_t} \tag{16}$$
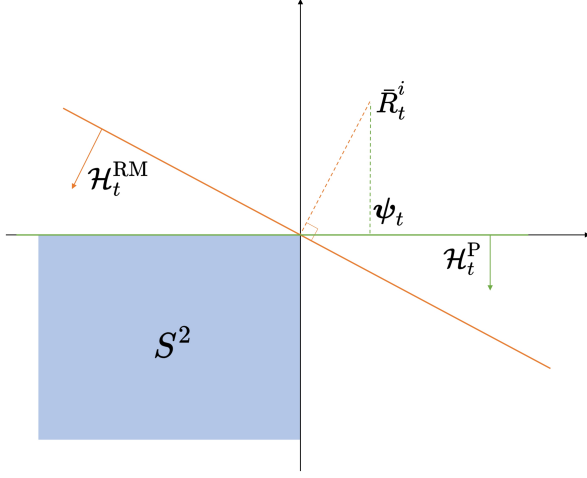
Figure 1: The difference between RM and FP (PCFR in normal form game) in a two-dimensional plane

vanilla CFR+ uses a weighted average strategy where iteration $T$ is weighted by $T$ rather than using a uniformly-weighted average strategy as in CFR.

**MCCFR** MCCFR also has many variants, but the most common is external sampling MCCFR(ES-MCCFR) because of its simplicity and powerful performance. In ES-MCCFR, some players are designated as traversers and others are samplers during an iteration. Traversers follow the CFR algorithm to update the regret and the average strategy of the experienced information set. On the rest of the sampler's nodes and the chance node, only one action is explored (sampled according to the player's strategy for that iteration on the information set), and the regret and the average strategy are not updated.

**CFR Pruning and CFR-BR** Pruning is a common optimization method in all tree search algorithms. The Alpha-Beta algorithm in the complete information extensive game is a pruned version for Min-Max algorithm. This method is also the key algorithm for Deep Blue's success (Hsu 2002). Naive pruning algorithm is the simplest pruning algorithm in CFR. In naive pruning, if all the players have no probability of reaching the current state $s$ ($\forall i, \pi_{\sigma_t}^{-i}(s) = 0$), the entire subtree at that state $s$ can be pruned for the current iteration without affecting the regret calculation. Existing research experiments have shown that only using the naive pruning algorithm can save more than 96% of the calculation time in some games (Lanctot et al. 2010).

CFR-BR (Johanson et al. 2012) is a variant of CFR in which one player takes an update of the CFR and the other player reacts optimally to the opponent's strategy in each iteration. In the game of perfect recall, calculating the best response to a fixed strategy is relatively simple computationally, just traversing the game tree once.

## Pure CFR Algorithm

### Fictitious Play achieves Blackwell approachability

We first prove that FP achieves Blackwell approachability in a two-player zero-sum game. Define $\bar{R}_t^{i,\max}$ be the maximum portion of vector $\bar{R}_t^{i,+}$. If $\bar{R}_t^{i,\max} \neq \mathbf{0}$, we find the point $\boldsymbol{\psi}_t = \bar{R}_t^i - \bar{R}_t^{i,\max}$ in $\mathbb{R}^{l^i}$ and let $\frac{\bar{R}_t^i - \boldsymbol{\psi}_t}{|\bar{R}_t^i - \boldsymbol{\psi}_t|}$ be the normal vector, then we can determine half-space by $\frac{\bar{R}_t^i - \boldsymbol{\psi}_t}{|\bar{R}_t^i - \boldsymbol{\psi}_t|}$ and $\boldsymbol{\psi}_t$:

$$\mathcal{H}_t^{\mathrm{P}} = \left\{ \boldsymbol{z} \in \mathbb{R}^{l^{-i}} : (\bar{R}_t^i - \boldsymbol{\psi}_t)^\top \boldsymbol{z} \leq (\bar{R}_t^i - \boldsymbol{\psi}_t)^\top \boldsymbol{\psi}_t \right\}$$
(17)

Because $\bar{R}_t^i - \boldsymbol{\psi}_t = \bar{R}_t^{i,\max}$, $(\bar{R}_t^i - \boldsymbol{\psi}_t)^\top \boldsymbol{\psi}_t = 0$, therefore:

$$\mathcal{H}_t^{\mathrm{P}} = \left\{ \boldsymbol{z} \in \mathbb{R}^{l^{-i}} : \left\langle \bar{R}_t^{i,\max}, \boldsymbol{z} \right\rangle \leq 0 \right\}$$
(18)

For any point $\boldsymbol{s}' \in S^i$ there is $\left\langle \bar{R}_t^{i,\max}, \boldsymbol{s}' \right\rangle \leq 0$. Then we need to prove that a forcing action for $\mathcal{H}_t^{\mathrm{P}}$ indeed exists. According to Definition 2, we need to find a $\sigma_{t+1}^{i*} \in \Sigma^i$ that achieves $R^i\left(\sigma_{t+1}^{i*}, \hat{\sigma}_{t+1}^{-i}\right) \in \mathcal{H}_{t+1}^{i,\mathrm{P}}$ for any $\hat{\sigma}_{t+1}^{-i} \in \Sigma^{-i}$. For simplicity, let $\boldsymbol{\ell} = \left[u^i\left(a_1, \sigma^{-i}\right), u^i\left(a_2, \sigma^{-i}\right), \dots\right]^\top \in \mathbb{R}^{l^i}$, we rewrite the regret vector as $R^i\left(\sigma_{t+1}^{i*}, \hat{\sigma}_{t+1}^{-i}\right) = \boldsymbol{\ell} - \left\langle \boldsymbol{\ell}, \sigma_{t+1}^{i*} \right\rangle \mathbf{1}$, we are looking for a $\sigma_{t+1}^{i*} \in \Sigma^i$ such that:

$$R^i\left(\sigma_{t+1}^{i*}, \hat{\sigma}_{t+1}^{-i}\right) \in \mathcal{H}_t^{\mathrm{P}}$$
$$\iff \left\langle \bar{R}_t^{i,\max}, \boldsymbol{\ell} - \left\langle \boldsymbol{\ell}, \sigma_{t+1}^{i*} \right\rangle \mathbf{1} \right\rangle \leq 0$$
$$\iff \left\langle \bar{R}_t^{i,\max}, \boldsymbol{\ell} \right\rangle - \left\langle \boldsymbol{\ell}, \sigma_{t+1}^{i*} \right\rangle \left\langle \bar{R}_t^{i,\max}, \mathbf{1} \right\rangle \leq 0$$
$$\iff \left\langle \bar{R}_t^{i,\max}, \boldsymbol{\ell} \right\rangle - \left\langle \boldsymbol{\ell}, \sigma_{t+1}^{i*} \right\rangle \left\| \bar{R}_t^{i,\max} \right\|_1 \leq 0 \quad (19)$$
$$\iff \left\langle \boldsymbol{\ell}, \frac{\bar{R}_t^{i,\max}}{\left\| \bar{R}_t^{i,\max} \right\|_1} \right\rangle - \left\langle \boldsymbol{\ell}, \sigma_{t+1}^{i*} \right\rangle \leq 0$$
$$\iff \left\langle \boldsymbol{\ell}, \frac{\left[\bar{R}_t^i\right]^{\max}}{\left\| \left[\bar{R}_t^i\right]^{\max} \right\|_1} - \sigma_{t+1}^{i*} \right\rangle \leq 0$$

We obtain the strategy $\sigma_{t+1}^{i*} = \frac{\bar{R}_t^{i,\max}}{\left\| \bar{R}_t^{i,\max} \right\|_1}$ that guarantees $\mathcal{H}_{t+1}^{\mathrm{P}}$ to be forceable half space. And the action with the highest regret value is actually the BR strategy (Brown 2020), so FP achieves Blackwell approachability. Figure 1 shows the difference in forceable half spaces for BR and RM strategies in the two-dimensional plane. According to Theorem 2, BR strategy is also a regret minimizer in normal form game.

### Algorithm Implementation Details of PCFR in Extensive Form Game

Since FP is a regret minimizer in normal form game, it means that replacing the RM strategy in the CFR with the BR strategy will not affect the convergence in extensive form game. So we propose a new algorithm called Pure CFR (PCFR).

First, it is no longer necessary to calculate the counterfactual regret, just define the immediate counterfactual value:

$$Q_{t,\text{imm}}^i(I,a) = Q_{t-1,\text{imm}}^i(I,a) + \pi_{\sigma_t}^{-i}(I)u^i\left(I, \sigma_t|_{I \to a}\right)$$
(20)

The difference between $Q_{t,\text{imm}}^i$ and $\bar{R}_{i,\text{imm}}^i$ is that $Q_{t,\text{imm}}^i$ does not need to subtract average payoff $u^i(I, \sigma_t)$, because it has no effect on finding the maximum value of $Q_{t,\text{imm}}^i$, and it will save a lot of computing time.

Second, the strategy in the next stage is a BR strategy rather than RM strategy:

$$\sigma_{t+1}^i = \arg\max_{a \in A} Q_{t,\text{imm}}^i(I,a)$$
(21)

Since $\arg\max$ will only take the action with the largest counterfactual value, the strategy obtained in each iteration is a pure strategy. This is why we named this algorithm Pure CFR. In addition, the time complexity of calculating the BR strategy is obviously lower than that of calculating the RM strategy.

Finally, since the update $\sigma_{t+1}^i$ is a pure strategy, $\pi_{\sigma_t}^{-i}(I)$ is either 0 or 1. It can be seen from the formula (20) that if $\pi_{\sigma_t}^{-i}(I) = 0$, it is unnecessary to calculate $Q_{t,\text{imm}}^i(I,a)$. Similarly, it can be seen from the formula (15) that if $\pi_{\sigma_t}^i(I) = 0$ it is unnecessary to update the historical strategy $\bar{\sigma}_t^i(I)$. These omitted steps will also greatly increase the efficiency of the algorithm. The pseudocode of the PCFR algorithm can be found in Appendix A.1.

## Pure CFR Algorithm Theoretical Analysis
### PCFR combined with CFR variants

The CFR variants, such as MCCFR, CFR+, and different average weighting schemes, are all applicable to the proposed PCFR (called PCFR-variants), because our algorithm meets Blackwell approachability. However, the theoretical convergence speed of the PCFR-variants can only be calculated in the worse cases. Thus, we use a series of experiments to explore the convergence speed of different PCFR-variants. According to the experiment results in Appendix B, the log-weighted PMCCFR (i.e., MCCFR + PCFR) converges the fastest. The pseudocode of the log-PMCCFR algorithm can be found in Appendix A.2.

### Time complexity reduction

Firstly, the core difference between PMCCFR and MCCFR is that PMCCFR directly takes the action with the largest regret value as the pure strategy for the next iteration, which greatly simplifies the algorithm. We compare the runtime complexity of MCCFR and PMCCFR in Appendix C.1. In theory, PMCCFR takes at most 2/9 of the original MCCFR's time when touching through the same number of nodes, and even taking into account the time of game simulation in training, the final PMCCFR only consumes about 1/3 to 1/4 of the time of MCCFR when passing through the same nodes. Since PMCCFR and MCCFR are met Blackwell approachability, their theoretical convergence speed is consistent, which makes the solution speed of PMCCFR 3–4 times faster than MCCFR in almost all problems.

Secondly, since the PMCCFR algorithm naturally rejects strictly dominated strategies, a new warm-start algorithm can be developed using this feature. The experiments in Appendix D show that too many dominated strategies are an important reason for the slow convergence rate of CFR+ in the early stage, and the convergence rate of CFR+ will be greatly accelerated after moving the dominated strategies. Correspondingly, the characteristic of PCFR is that if a strategy has not been selected or has only been selected a few times in a long period of iterations, then this strategy has a high probability of being a dominated strategy. We combine the advantages of PCFR and CFR+ to develop a new warm-start method–eliminating dominated strategies. In the early $T_e$ iteration training, we use PMCCFR to obtain an approximate equilibrium strategy $\sigma_e$. If $\exists a, \sigma_e(a) < \xi$ ($\xi$ is a small probability), then the action $a$ has a large probability of being a strictly dominated strategy. We remove these strategies, then use the CFR+ algorithm to train the game after eliminating dominated strategy. Experiments show that when the parameters $T_e$ and $\xi$ are selected appropriately, this method will greatly accelerate the convergence speed.

### Space complexity reduction

PMCCFR will prune most of the actions, and we show in Appendix C.2 that PMCCFR only touch $\mathcal{O}(\sqrt{|\mathcal{S}|})$ nodes in one iteration. This makes PMCCFR applicable to the early approximate equilibrium solving of all game problems without any hardware threshold. Taking advantage of this feature, PMCCFR can obtain many approximate equilibriums by parallel computing, and then combine these approximate equilibrium into a stronger strategy. Since there are still many nodes in one iteration, this parallel idea is difficult to achieve in CFR or MCCFR in large games.

## Experiments Evaluation
### Description of the game

We have adopted matrix games, Kuhn-extension poker, Leduc-extension poker, princess and monster (Lanctot et al. 2010) to compare the abilities of different algorithms. The Kuhn-extension poker is based on the vanilla Kuhn poker (Kuhn 1950), increasing the number of cards to $x$ (vanilla Kuhn poker has 3 cards); increasing the bet action types to $y$ (only 1 in vanilla Kuhn poker); increasing the number of raising times to $z$ times (vanilla Kuhn poker has only 1 time); The improvement of Leduc-extension is similar to the improvement of Kuhn-extension based on Leduc poker (Shi and Littman 2002). These improvements allow us to change the scale of the game very easily. Princess and monster is a classic pursuit problem. For a detailed description of these games, and more experimental results, see Appendix E.

### Experimental settings

The CPU used in our experiment is AMD 3990WX, and the memory is 128GB. We run a set of experiments to compare PMCCFR with CFR, ES-MCCFR, and CFR+.

In order to reflect the fluctuation range of convergence of different algorithms under different conditions, all algo-

rithms adopt a random distribution as the initial strategy and update the strategies of all players synchronously in one iteration. The weight settings of the comparative experiments are consistent with previous classic papers (Lanctot et al. 2010; Brown 2020). ES-MCCFR and CFR both adopt square weights, the time weights of CFR+ are set to liner, while PCFR and PMCCFR adopt log weights. In the engineering implementation, if the probability of reaching a certain node during iteration is less than $10^{-20}$, it will be directly pruned. In MCCFR, if $R_t^{i,\max} = \mathbf{0}$, we directly set the strategy of the next stage to a randomly pure strategy.

In Appendix F, we compare the advantages and disadvantages of the time, number of iterations, and the number of passing nodes as a measurement index. Finally, we will measure the convergence of the algorithm from the two perspectives of passing nodes and time.

## Experimental results

**Convergence**  As shown in the Figure 2, when the number of nodes touched is used as the horizontal axis, the convergence curve of the PMCCFR method is slightly better than that of MCCFR. We speculate that this is because both algorithms satisfy Blackwell approachability, so their convergence rates are similar. However, the advantage of PMCCFR is that when passing the same number of nodes, the computation time is reduced to about 2/3 of MCCFR due to the engineering implementation trick mentioned above. Therefore, when reaching the same exploitability, PMCCFR will save 3/4~2/3 of the time compared with MCCFR.

The Convergence speed in 14-Card 1-Action 1-Len Leduc shows that although MCCFR will surpass CFR+ when the passing nodes is used as the horizontal axis, CFR+ will surpass MCCFR when time is used as the horizontal axis. This is because MCCFR calculates the RM strategy every time it passes through a node, while CFR+ only calculates the RM strategy once after passing through all nodes in an information set. Therefore, when passing through the same number of nodes, the time of MCCFR calculate RM strategy will definitely exceed CFR+, which also reflects the necessity of taking time as the horizontal axis.

**Comparison of time and space required for post-pruning game traversal**  CFR+ needs to traverse the entire game tree for each iteration, which makes it extremely difficult to use the CFR+ algorithm in large-scale games. Correspondingly, if a small number of nodes are passed in one iteration, it can not only allow the algorithm to run on lower hardware configurations, but also does not need to consider any engineering problems such as memory allocation and recycling, which will bring great convenience to engineering implementation. Table 1 compares the number of nodes touched in the first five iterations of different games (averaged over 30 random samples). As shown in the Table 1, in the early iterations of training, CFR+/CFR needs to traverse about half of the nodes of the entire game tree, while the PMCCFR traverses very few nodes. In the 3-Card 20-Action 7-Len Kuhn with 4.97 million nodes, PMCCFR will only pass 855 nodes at one iteration. In this game, the sampling algorithm MCCFR will pass through 5528.12 nodes, but the full traver-

| Game name | Game tree's node | Travesal or MC | Algorithm | Average number of nodes touched per itration | The proportion of nodes touched per iteration (%) |
|---|---|---|---|---|---|
| 5 Card 5 Bet Action 1 Len Leduc | 245336 | Traversal | CFR+ | 131615.8 | 53.647 |
| | | | CFR | 110428.3 | 45.011 |
| | | | PCFR | 11027.4 | 4.494 |
| | | MC | PMCCFR | **110.4** | **0.045** |
| | | | MCCFR | 221.3 | 0.09 |
| 3 Card 20 Bet Action 7 Len Kuhn | 4967263 | Traversal | CFR+ | 2774844.0 | 55.862 |
| | | | CFR | 2198996.4 | 44.269 |
| | | | PCFR | 5413.0 | 0.108 |
| | | MC | PMCCFR | **855.8** | **0.017** |
| | | | MCCFR | 5528.1 | 0.111 |

Table 1: The number of nodes that different algorithms pass through in one iteration

sal algorithm PCFR will only pass through 5413.04 nodes. This means that as long as the training time is sufficient, PMCCFR can be directly applied to the solution to large-scale game problems without any memory recycling. In addition, previous experience has shown that algorithms that pass fewer nodes in one iteration often converge faster in training large-scale games, which also explains why PMCCFR is slightly faster than MCCFR when measuring convergence speed by the number of nodes passed.

**Warm Start PMCCFR**  We described a new warm start idea–eliminate dominated strategies. First, an approximate equilibrium $\sigma_e$ is generated using PMCCFR in the early stages of training. Then, $a$ can be considered as a dominated strategy if the probability of choosing action $a$ is less than $\xi$ ($\sigma_e(a) < \xi$), which means that eliminating action $a$ will not affect the strategy convergence to NE. Finally, using CFR+ to solve the game equilibrium after the dominated strategies is eliminated.

We use matrix games to test the effect of warm start. First, using a Gaussian distribution with mean 0 and variance 1 to generate a $100 \times 100$ matrix, adding 2 to the first $j$ strategies of player1 in the matrix, and subtracting 2 to the first $j$ strategies of player2 (this means that the last $100 - j$ strategies are highly likely to be dominated strategies). PMCCFR,CFR+ and CFR+ with warm start, are used to calculate the equilibrium of the game respectively, where warm start has two modes of MCCFR and PMCCFR. As shown in the Figure 3, when the proportion of the dominated strategies is low, the convergence speed of the warm start method based on PMCCFR surpasses the original CFR+ by two orders of magnitude, and PMCCFR only needs less than 5% of the CFR+ training time to obtain the warm start strategy $\sigma_e$. Although theoretically, the MCCFR algorithm will also eliminate the dominated strategies after long-term training, but experiments show that MCCFR is far less effective than PMCCFR as a warm start for eliminating the dominated strategies. Because in RM, as long as a strategy's payoff is greater than the average, then this strategy has the probability of being adopted. However, the payoff of dominated strategies is likely to be greater than the average, so the RM algorithm is difficult to use as a warm start with eliminating the dominated strategies.
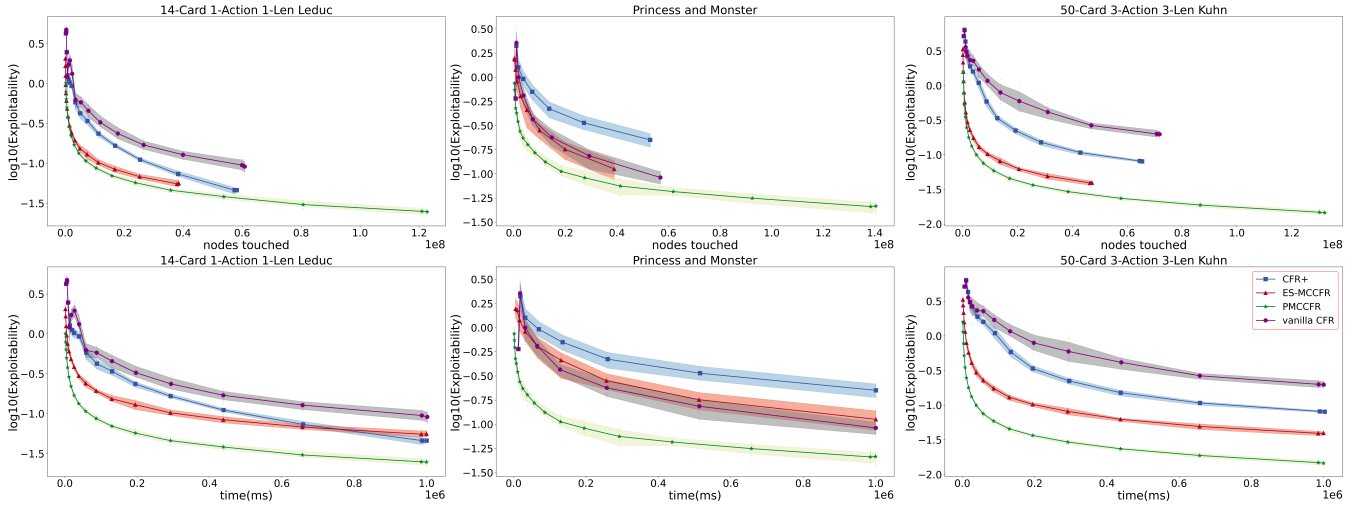
Figure 2: Convergence speed in Leduc-extension, Princess and monster, Kuhn-extension. The graph in the first row takes the number of passing nodes as the horizontal axis, and the graph in the second row takes the running time of the algorithm as the horizontal axis. The training timing of all experiments is fixed at 1000s. Each experiment has an average of 30 rounds, and the light range is the 90% confidence interval.
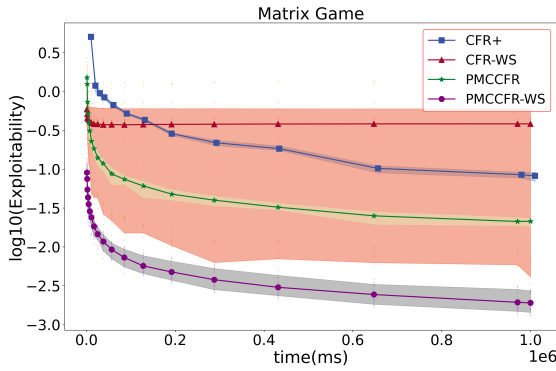


Figure 3: Convergence in Martrix Game with $\xi = \frac{10}{T}$, $j = 5$

## Conclusion

This paper proposes a novel method for solving incomplete information zero-sum games——Pure CFR. We first demonstrate that PCFR still has the same convergence properties and a similar convergence rate as the vanilla CFR, and can be freely combined with previous CFR variants. Due to adopting the BR strategy instead of the RM strategy, PMC-CFR achieved a speed increase of $3 \sim 4$ times compared to vanilla MCCFR in the final experiment. At the same time, we found that PCFR completely avoids strictly dominated strategies and proposed a warm start method based on this. This method can achieve two orders of magnitude improvement in games with a high proportion of dominated strategies.

In the future research, we will verify the feasibility of the method on a larger scale and combine it with other CFR variants, such as Alternate Update, Discount CFR (Brown and

Sandholm 2019a), Lazy CFR (Zhou et al. 2018), Greedy RM (Zhang, Lerer, and Brown 2022), and predictive version (Farina, Kroer, and Sandholm 2021) etc. In addition, since fewer nodes pass through an iteration, another focus is to explore the ability of this method in parallel computing and in combination with deep networks.

In the warm start of PCFR, as the proportion of dominant strategies increases, the convergence acceleration brought by this warm start decreases, and the convergence instability increases (unable to converge to the correct NE). This is because once any dominant strategy is eliminated, it will not be able to converge to NE in the end. In subsequent research, the warm start of eliminating dominated strategies needs to solve two problems: when to switch from PMCCFR to CFR+ given the training time? What is the appropriate probability $\xi$ of eliminating an action when switching from PCFR to CFR+? These issues remain to be further studied.

## References

Abernethy, J.; Bartlett, P.; and Hazan, E. 2011. Blackwell Approachability and No-Regret Learning are Equivalent.

Berger, U. 2007. Brown's original fictitious play. *Journal of Economic Theory*, 135(1): 572–578.

Bowling, M.; Burch, N.; Johanson, M.; and Tammelin, O. 2015. Heads-up limit hold'em poker is solved. *Science*, 347(6218): 145–149.

Brown, G. W. 1951. Iterative solution of games by fictitious play. *activity analysis of production and allocation*.

Brown, N. 2020. Equilibrium finding for large adversarial imperfect-information games. *PhD thesis*.

Brown, N.; and Sandholm, T. 2019a. Solving Imperfect-Information Games via Discounted Regret Minimization.

*Proceedings of the AAAI Conference on Artificial Intelligence*, 33: 1829–1836.

Brown, N.; and Sandholm, T. 2019b. Superhuman AI for multiplayer poker. *Science*, 365(6456): eaay2400.

Farina, G.; Grand-Clément, J.; Kroer, C.; Lee, C.-W.; and Luo, H. 2023. Regret Matching+: (In)Stability and Fast Convergence in Games. arXiv:2305.14709.

Farina, G.; Kroer, C.; and Sandholm, T. 2021. Faster game solving via predictive blackwell approachability: Connecting regret matching and mirror descent. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 5363–5371.

Fudenberg, D.; and Levine, D. K. 1998. The Theory of Learning in Games. *MIT Press Books*, 1.

Hart, S.; and Mas-Colell, A. 2000. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68(5): 1127–1150.

Heinrich, J.; Lanctot, M.; and Silver, D. 2015. Fictitious Self-Play in Extensive-Form Games. In *International Conference on Machine Learning*.

Hendon, E.; Jacobsen, H. J.; and Sloth, B. 1996. Fictitious play in extensive form games. *Games and Economic Behavior*, 15(2): 177–202.

Hsu, H. F. H. 2002. Deep Blue. *Artificial Intelligence*.

Johanson, M. 2013. Measuring the Size of Large No-Limit Poker Games. *Computer Science*.

Johanson, M.; Bard, N.; Burch, N.; and Bowling, M. 2012. Finding optimal abstract strategies in extensive-form games. In *National Conference on Artificial Intelligence*.

Kuhn, H. W. 1950. Simplified two-person poker. *Contributions to the Theory of Games*.

Lanctot, M.; Waugh, K.; Zinkevich, M.; and Bowling, M. H. 2010. Monte Carlo Sampling for Regret Minimization in Extensive Games. In *Advances in Neural Information Processing Systems 22*.

Leslie, D. S.; and Collins, E. J. 2006. Generalised weakened fictitious play. *Games and Economic Behavior*, 56(2): 285–298.

Li, J.; Koyamada, S.; Ye, Q.; Liu, G.; and Hon, H. W. 2020. Suphx: Mastering Mahjong with Deep Reinforcement Learning.

Moravík, M.; Schmid, M.; Burch, N.; Lis, V.; and Bowling, M. 2017. DeepStack: Expert-Level Artificial Intelligence in No-Limit Poker. *Science*, 356(6337): 508.

Shi, J.; and Littman, M. L. 2002. Abstraction Methods for Game Theoretic Poker. In *International Conference on Computers and Games*.

Silver, D.; Hubert, T.; Schrittwieser, J.; Antonoglou, I.; Lai, M.; Guez, A.; Lanctot, M.; Sifre, L.; Kumaran, D.; Graepel, T.; Lillicrap, T.; Simonyan, K.; and Hassabis, D. 2018. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science*, 362(6419): 1140–1144.

Tammelin, O. 2014. Solving Large Imperfect Information Games Using CFR+. *Eprint Arxiv*.

Zha, D.; Lai, K. H.; Cao, Y.; Huang, S.; Wei, R.; Guo, J.; and Hu, X. 2019. RLCard: A Toolkit for Reinforcement Learning in Card Games.

Zha, D.; Xie, J.; Ma, W.; Zhang, S.; and Liu, J. 2021. DouZero: Mastering DouDizhu with Self-Play Deep Reinforcement Learning.

Zhang, H.; Lerer, A.; and Brown, N. 2022. Equilibrium Finding in Normal-Form Games Via Greedy Regret Minimization.

Zhou, Y.; Ren, T.; Li, J.; Yan, D.; and Zhu, J. 2018. Lazy-CFR: fast and near optimal regret minimization for extensive games with imperfect information.

Zinkevich, M.; Johanson, M.; Bowling, M.; and Piccione, C. 2007. Regret minimization in games with incomplete information. *Oldbooks.nips.cc*, 20: 1729–1736.