

Double Neural Counterfactual Regret Minimization

Hui Li¹, Kailiang Hu¹, Zhibang Ge¹, Tao Jiang¹, Yuan Qi¹, and Le Song^{1,2}

Ant Financial¹

Georgia Institute of Technology²

{ken.lh, hkl163251, zhibang.zg, lvshan.jt, yuan.qi, le.song}@antfin.com
lsong@cc.gatech.edu

Abstract

Counterfactual regret minimization (CRF) is a fundamental and effective technique for solving imperfect information games. However, the original CRF algorithm only works for discrete state and action spaces, and the resulting strategy is maintained as a tabular representation. Such tabular representation limits the method from being directly applied to large games and continuing to improve from a poor strategy profile. In this paper, we propose a double neural representation for the Imperfect Information Games, where one neural network represents the cumulative regret, and the other represents the average strategy. Furthermore, we adopt the counterfactual regret minimization algorithm to optimize this double neural representation. To make neural learning efficient, we also developed several novel techniques including a robust sampling method, mini-batch Monte Carlo counterfactual regret minimization (MCCFR) and Monte Carlo counterfactual regret minimization plus (MCCFR+) which may be of independent interests. Experimentally, we demonstrate that the proposed double neural algorithm converges significantly better than the reinforcement learning counterpart.¹

Introduction

In Imperfect Information Games (IIG), a player only has partial access to the knowledge of her opponents before making a decision. This is similar to the real world scenarios, such as trading, traffic routing, and public auction. Thus designing methods for solving IIG is of great economic and social benefits.

Nash equilibrium is a typical solution concept for a two-player extensive-form game. CFR (Zinkevich et al. 2007) is an efficient algorithm that approximate Nash equilibrium in large games. CFR try to minimize overall counterfactual regret and prove that the average of the strategies in all iterations would converge to a Nash equilibrium. However, the original CFR only works for discrete state and action spaces, and the resulting strategy is maintained as a tabular representation. Such tabular representation limits the method from being directly applied to large games and continuing to improve if starting from a poor strategy profile.

To alleviate CFR’s large memory requirement in large games such as heads-up no-limit Texas Hold’em, (Moravk et al. 2017) proposed a seminal approach called DeepStack which uses fully connected neural networks to represent players’ counterfactual values and obtain a strategy online as requested. However, the strategy is still represented as a tabular form and the quality of this solution depends a lot on the initial quality of the counterfactual network. Furthermore, the counterfactual network is estimated separately, and it is not easy to continue improving both counterfactual network and the tabular strategy profile in an end-to-end optimization framework. (Heinrich, Lanctot, and Silver 2015) or (Heinrich and Silver 2016) proposed end-to-end fictitious self-play approaches (XFP and NFSP respectively) to learn the approximate Nash equilibrium with deep reinforcement learning. In a fictitious play model, strategies are represented as neural networks and the strategies are updated by selecting the best responses to their opponents’ average strategies. This approach is advantageous in the sense that the approach does not rely on abstracting the game, and in theory, the strategy should continually improve as the algorithm iterates more steps. However, these methods do not explicitly take into account the hidden information in a game, and in experiments for games such as Leduc Hold’em, these methods converge slower than tabular based counterfactual regret minimization algorithms. (Vaughan et al. 2015) used hand-crafted features of the information sets to estimate the counterfactual regret. However, it need traversals of the full game tree which is infeasible in large games.

Thus it remains an open question whether the purely neural-based end-to-end approach can achieve comparable performance to tabular based CFR approach. In the paper, we partially resolve this open question by designing a double neural counterfactual regret minimization algorithm which can match the performance of tabular based counterfactual regret minimization algorithm. We employed two neural networks, one for the cumulative regret, and the other for the average strategy. We show that careful algorithm design allows these two networks to track the cumulative regret and average strategy respectively, resulting in a converging neural strategy. Furthermore, in order to improve the convergence of the neural algorithm, we also developed a new sampling technique which has lower variance than the outcome sampling, while being more memory efficient than the exter-

nal sampling. In experiments with Leduc Hold'em and One-card poker, we showed that the proposed double neural algorithm can converge to comparable results produced by its tabular counterpart while performing much better than deep reinforcement learning method. The current results open up the possibility for a purely neural approach to directly solve large IIG.

Background

Representation of Extensive-Form Game

We define the components of an extensive-form game following (Osborne and Ariel 1994). A finite set $N = \{0, 1, \dots, n-1\}$ of **players**. Define h_i^v as the **hidden variable** of player i in IIG. H refers to a finite set of histories. Each member $h = (h_i^v)_{i=0,1,\dots,n-1}(a_l)_{l=0,\dots,L-1} = h_0^v h_1^v \dots h_{n-1}^v a_0 a_1 \dots a_{L-1}$ of H denotes a possible **history** (or state), which consists of each player's hidden variable and L actions taken by players including chance. For player i , h also can be denoted as $h_i^v h_{-i}^v a_0 a_1 \dots a_{L-1}$, where h_{-i}^v refers to the opponent's hidden variables. The empty sequence \emptyset is a member of H . $h_j \sqsubseteq h$ denotes h_j is a prefix of h . $Z \subseteq H$ denotes the terminal histories and any member $z \in Z$ is not a prefix of any other sequences. $A(h) = \{a : ha \in H\}$ is the set of available actions in non-terminal history $h \in H \setminus Z$. A **player function** $P(h)$ returns the player who acts at history $h \in N \cup \{c\}$, where c denotes the chance, which usually is -1. $P(h)$ is the player who takes an action after history h . \mathcal{I}_i of a history $\{h \in H : P(h) = i\}$ is an **information partition** of player i . A set $I_i \in \mathcal{I}_i$ is an **information set** of player i and $I_i(h)$ refers to information set I_i at state h . Generally, I_i indicates a sequence in IIG, *i.e.*, $h_i^v a_0 a_1 \dots a_{L-1}$. For $I_i \in \mathcal{I}_i$ we denote by $A(I_i)$ the set $A(h)$ and by $P(I_i)$ the player $P(h)$ for any $h \in I_i$. For each player $i \in N$ a utility function $u_i(z)$ defines the payoff of the terminal state z .

Strategy and Nash equilibrium

The strategy in an extensive-form game contains the following components. A **strategy profile** $\sigma = \{\sigma_i | \sigma_i \in \Sigma_i, i \in N\}$ is a collection of strategies for all players, where Σ_i is the set of all possible strategies for player i . σ_{-i} refers to all strategies in σ except σ_i . For play $i \in N$ the **strategy** $\sigma_i(I_i)$ is a function, which assigns an action distribution over $A(I_i)$ to information set I_i . $\sigma_i(a|h)$ denotes the probability of action a taken by player i at $N \cup \{c\}$ at state h . In IIG, $\forall h_1, h_2 \in I_i$, we have $I_i = I_i(h_1) = I_i(h_2)$, $\sigma_i(I_i) = \sigma_i(h_1) = \sigma_i(h_2)$, $\sigma_i(a|I_i) = \sigma_i(a|h_1) = \sigma_i(a|h_2)$. For iterative method such as CFR, σ^t refers to the strategy profile at t -th iteration. The **state reach probability** of history h is denoted by $\pi^\sigma(h)$ if players take actions according to σ . For an empty sequence $\pi^\sigma(\emptyset) = 1$. The reach probability can be decomposed into $\pi^\sigma(h) = \pi_i^\sigma(h) \pi_{-i}^\sigma(h)$ according to each player's contribution, where $\pi_i^\sigma(h) = \prod_{h' a \sqsubseteq h, P(h')=P(h)} \sigma_i(a|h')$ and $\pi_{-i}^\sigma(h) = \prod_{h' a \sqsubseteq h, P(h') \neq P(h)} \sigma_{-i}(a|h')$. The **information set reach probability** of I_i is defined as $\pi^\sigma(I_i) = \sum_{h \in I_i} \pi^\sigma(h)$. If $h' \sqsubseteq h$, the **interval state reach probability** from state h' to h is defined as $\pi^\sigma(h', h)$, then we have

$\pi^\sigma(h', h) = \pi^\sigma(h) / \pi^\sigma(h')$. $\pi_i^\sigma(I_i)$, $\pi_{-i}^\sigma(I_i)$, $\pi_i^\sigma(h', h)$, and $\pi_{-i}^\sigma(h', h)$ are defined similarly.

Counterfactual Regret Minimization

In large and zero-sum IIG, CFR is proved to be an efficient method to compute Nash equilibrium (Zinkevich et al. 2007), (Brown and Sandholm 2017) or (Moravk et al. 2017). We present some key ideas of this method as follows.

Lemma 1: The state reach probability of one player is proportional to posterior probability of the opponent's hidden variable, *i.e.*, $p(h_{-i}^v | I_i) \propto \pi_{-i}^\sigma(h)$, where h_{-i}^v and I_i indicate a particular h . (We provide the proof in (Li et al. 2018).)

For player i and strategy profile σ , the **counterfactual value (CFV)** $v_i^\sigma(h)$ at state h is define as

$$v_i^\sigma(h) = \sum_{h \sqsubseteq z, z \in Z} \pi_i^\sigma(h, z) u_i'(z). \quad (1)$$

where $u_i'(z) = \pi_{-i}^\sigma(z) u_i(z)$ is the expected reward of player i with respect to the approximated posterior distribution of the opponent's hidden variable. The **action counterfactual value** of taking action a is $v_i^\sigma(a|h) = v_i^\sigma(ha)$ and the regret of taking this action is $r_i^\sigma(a|h) = v_i^\sigma(a|h) - v_i^\sigma(h)$. Similarly, the CFV of information set I_i is $v_i^\sigma(I_i) = \sum_{h \in I_i} v_i^\sigma(h)$ and the regret is $r_i^\sigma(a|I_i) = \sum_{z \in Z, ha \sqsubseteq z, h \in I_i} \pi_i^\sigma(ha, z) u_i'(z) - \sum_{z \in Z, h \sqsubseteq z, h \in I_i} \pi_i^\sigma(h, z) u_i'(z)$. Then the **cumulative regret** of action a after T iterations is

$$R_i^T(a|I_i) = R_i^{T-1}(a|I_i) + r_i^{\sigma^T}(a|I_i). \quad (2)$$

where $R_i^0(a|I_i) = 0$. Define $R_i^{T,+}(a|I_i) = \max(R_i^T(a|I_i), 0)$, the **current strategy** at $T+1$ iteration will be updated by

$$\sigma_i^{T+1}(a|I_i) = \begin{cases} \frac{R_i^{T,+}(a|I_i)}{\sum_{a \in A(I_i)} R_i^{T,+}(a|I_i)} & \text{if } \sum_{a \in A(I_i)} R_i^{T,+}(a|I_i) > 0 \\ \frac{1}{|A(I_i)|} & \text{otherwise.} \end{cases} \quad (3)$$

The **average strategy** $\bar{\sigma}_i^T$ after T iterations is defined as:

$$\bar{\sigma}_i^T(a|I_i) = \frac{\sum_{t=1}^T \pi_i^{\sigma^t}(I_i) \sigma_i^t(a|I_i)}{\sum_{t=1}^T \pi_i^{\sigma^t}(I_i)}. \quad (4)$$

where $\pi_i^{\sigma^t}(I_i)$ denotes the information set reach probability of I_i at t -th iteration and is used to weight the corresponding current strategy $\sigma_i^t(a|I_i)$. Define $s_i^t(a|I_i) = \pi_i^{\sigma^t}(I_i) \sigma_i^t(a|I_i)$ as the additional numerator in iteration t , then the cumulative numerator can be defined as $S^T(a|I_i) = S^{T-1}(a|I_i) + s_i^T(a|I_i)$, where $S^0(a|I_i) = 0$.

Monte Carlo CFR

Vanilla CFR needs to traverse the entire game tree which is infeasible in large game. (Lanctot et al. 2009) proposed a sample-based algorithms called Monte Carlo counterfactual regret minimization(MCCFR) which traverse part of game tree. Define $\mathcal{Q} = \{Q_1, Q_2, \dots, Q_m\}$, where $Q_j \in Z$ is a **block** of sampling terminal histories in each iteration. Define q_{Q_j} as the probability of considering block Q_j , where

$\sum_{j=1}^m q_{Q_j} = 1$. Define $q(z) = \sum_{j:z \in Q_j} q_{Q_j}$ as the probability of considering a particular terminal history z . For information set I_i , a sample estimate of **counterfactual value** is $\tilde{v}_i^\sigma(I_i|Q_j) = \sum_{h \in I_i, z \in Q_j, h \sqsubseteq z} \frac{1}{q(z)} \pi_{-i}^\sigma(z) \pi_i^\sigma(h, z) u_i(z)$.

Lemma 2: $E_{j \sim q_{Q_j}} [\tilde{v}_i^\sigma(I_i|Q_j)] = v_i^\sigma(I_i)$ (see the proof of Lemma 1 in (Lanctot et al. 2009).)

Define $\sigma^{r,s}$ as **sampling strategy profile**, where $\sigma_i^{r,s}$ is the sampling strategy for player i and $\sigma_{-i}^{r,s}$ are the sampling strategies for players expect i . Particularly, for both external sampling and outcome sampling proposed by (Lanctot et al. 2009), $\sigma_{-i}^{r,s} = \sigma_{-i}$. The regret of the sampled action $a \in A(I_i)$ is defined as

$$\begin{aligned} \tilde{r}_i^\sigma((a|I_i)|Q_j) &= \sum_{z \in Q_j, ha \sqsubseteq z, h \in I_i} \pi_i^\sigma(ha, z) u_i^{r,s}(z) \\ &\quad - \sum_{z \in Q_j, h \sqsubseteq z, h \in I_i} \pi_i^\sigma(h, z) u_i^{r,s}(z) \end{aligned} \quad (5)$$

where $u_i^{r,s}(z) = u_i(z)/\pi_i^{\sigma^{r,s}}(z)$ is a new utility weighted by $1/\pi_i^{\sigma^{r,s}}(z)$. The sample estimate for **cumulative regret** of action a after T iterations is $\tilde{R}_i^T((a|I_i)|Q_j) = \tilde{R}_i^{T-1}((a|I_i)|Q_j) + \tilde{r}_i^{\sigma^T}((a|I_i)|Q_j)$ with $\tilde{R}_i^0((a|I_i)|Q_j) = 0$.

Double Neural Counterfactual Regret Minimization

In this section, we will explain our double neural CFR algorithm, where we employ two neural networks, one for the cumulative regret, and the other for the average strategy.

As shown in Figure 1 (A), standard CFR-family methods such as CFR (Zinkevich et al. 2007), outcome-sampling MCCFR, external sampling MCCFR (Lanctot et al. 2009), and CFR+ (Tammelin 2014), (Johanson et al. 2012), (Burch et al. 2012) need to use two large tabular-based memories \mathcal{M}_R and \mathcal{M}_S to record the cumulative regret and average strategy for all information sets. Such tabular representation makes these methods difficult to apply to large extensive-form games with limited time and space (Burch 2017).

In contrast, we will use two deep neural networks to compute approximate Nash equilibrium of IIG as shown in Figure 1 (B). Different from NFSP, our method is based on the theory of CFR, where the first network is used to learn the cumulative regret and the other is to learn the cumulative numerator of the average strategy profile. With the help of these two networks, we do not need to use a large memory

to save the key information of the entire game tree. In practice, the proposed double neural method can achieve a lower exploitability with fewer iterations than NFSP. In addition, we present experimentally that our double neural CFR can also continually improve after initialization from a poor tabular strategy.

Overall Framework

An algorithm in the CFR framework needs to be able to answer two queries:

1. what is the current strategy $\sigma^{t+1}(a|I_i)$ for iteration $t+1$;
2. and what is the average strategy $\bar{\sigma}_i^t(a|I_i)$ after t iterations; $\forall i \in N, \forall I_i \in \mathcal{I}_i, \forall a \in A(I_i), \forall t \in [1, T]$. Thus, our neural networks are designed to address the needs for these two queries respectively.

For the first query. According to Eq. (3), current strategy $\sigma^{t+1}(a|I_i)$ is computed by the cumulative regret $R^t(a|I_i)$. Given information set I_i and action a , we design a neural network **RegretSumNetwork(RSN)** $\mathcal{R}(a, I_i|\theta_{\mathcal{R}}^t)$ to learn $R^t(a|I_i)$, where $\theta_{\mathcal{R}}^t$ is the parameter in the network at t -th iteration. As shown Figure 1 (b), define memory $\mathcal{M}_R = \{(I_i, \tilde{r}_i^{\sigma^t}((a|I_i)|Q_j)) | \forall i \in N, \forall a \in A(I_i), h \in I_i, h \sqsubseteq z, z \in Q_j\}$. Each member of \mathcal{M}_R is the visited information set I_i and the corresponding regret $\tilde{r}_i^{\sigma^t}((a|I_i)|Q_j)$, where Q_j is the sampled block in t -th iteration. According to Eq. (2), we can estimate $\mathcal{R}(a, I_i|\theta_{\mathcal{R}}^{t+1})$ using the following optimization:

$$\begin{aligned} \theta_{\mathcal{R}}^{t+1} \leftarrow \operatorname{argmin}_{\theta_{\mathcal{R}}^{t+1}} \sum_{(I_i, \tilde{r}_i^{\sigma^t}((a|I_i)|Q_j)) \in \mathcal{M}_R} \left(\mathcal{R}(a, I_i|\theta_{\mathcal{R}}^t) \right. \\ \left. + \tilde{r}_i^{\sigma^t}((a|I_i)|Q_j) - \mathcal{R}(a, I_i|\theta_{\mathcal{R}}^{t+1}) \right)^2 \end{aligned} \quad (6)$$

For the second query. According to Eq. (4), the approximate Nash equilibrium is the weighted average of all previous strategies over T iterations. We only need to track the numerator in Eq. (4) since the denominator is used to normalize the summation. Similar to the cumulative regret, we employ another deep neural network **AvgStrategyNetwork(ASN)** to learn the cumulative numerator of the average strategy. Define $\mathcal{M}_S = \{(I_i, \pi_i^{\sigma^t}(I_i)\sigma_i^t(a|I_i)) | \forall i \in N, \forall a \in A(I_i), h \in I_i, h \sqsubseteq z, z \in Q_j\}$. Each member of \mathcal{M}_S is the visited information set I_i and the value of $\pi_i^{\sigma^t}(I_i)\sigma_i^t(a|I_i)$, where Q_j is the sampled block in t -th iteration. Then the parameter θ_S^{t+1} can estimated by the fol-

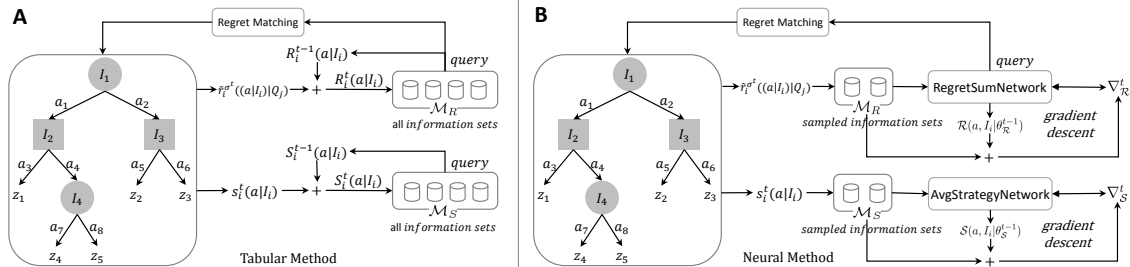


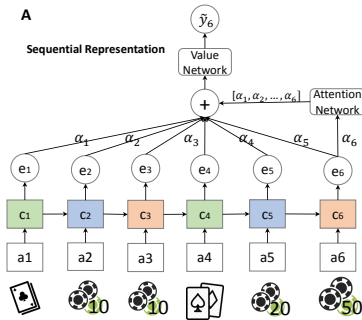
Figure 1: (A) tabular based CRF, and (B) our double neural based CRF framework.

lowing optimization:

$$\theta_S^{t+1} \leftarrow \underset{\theta_S^{t+1}}{\operatorname{argmin}} \sum_{(I_i, s_i^t(a|I_i)) \in \mathcal{M}_S} \left(\mathcal{S}(a, I_i | \theta_S^t) + s_i^t(a|I_i) - \mathcal{S}(a, I_i | \theta_S^{t+1}) \right)^2. \quad (7)$$

Recurrent Neural Network Representation

In order to define our \mathcal{R} and \mathcal{S} network, we need to represent the information set $I_i \in \mathcal{I}$ in extensive-form games. In such games, players take action in alternating fashion and each player makes a decision according to the observed history. In this paper, we model the behavior sequence as a recurrent neural network and each action in the sequence corresponds to a cell in RNN. Figure 2 (A) provides an illustration of the proposed deep sequential neural network representation for information sets. Here we will use LSTM for the representation. Furthermore, different position in the sequence may contribute differently to the decision making, we will add an attention mechanism (Desimone and Duncan 1995), (Cho, Courville, and Bengio 2015) to the LSTM architecture to enhance the representation. For example, the player may need to take a more aggressive strategy after beneficial public cards are revealed. Thus the information, after the public cards are revealed may be more important. More specifically, for l -th cell, define x_l as the input vector (which can be either player or chance actions), e_l as the hidden layer embedding, ϕ_* as a general nonlinear function. Each action is represented by a LSTM cell, which has the ability to remove or add information to the cell state with three different gates. Define the notation \cdot as element-wise product. The first **forgetting gate** layer is defined as $g_l^f = \phi_f(w^f[x_l, e_{l-1}])$, where $[x_l, e_{l-1}]$ denotes the concatenation of x_l and e_{l-1} . The second **input gate** layer decides which values to update and is defined as $g_l^i = \phi_i(w^i[x_l, e_{l-1}])$. A nonlinear layer output a vector of new candidate values $\tilde{C}_l = \phi_c(w^l[x_l, e_{l-1}])$ to decide what can be added to the state. After the forgetting gate and the input gate, the new cell state is updated by $C_l = g_l^f \cdot C_{l-1} + g_l^i \cdot \tilde{C}_l$. The third **output gate** is defined as $g_l^o = \phi_o(w^o[x_l, e_{l-1}])$. Finally, the updated hidden embedding is $e_l = g_l^o \cdot \phi_e(C_l)$. As shown in Figure 2 (A), for each LSTM cell j , the vector of attention weight is learned



by an **attention network**. Each member in this vector is a scalar $\alpha_j = \phi_a(w^a e_j)$. The attention embedding of l -th cell is then defined as $e_l^a = \sum_{j=1}^l \alpha_j \cdot e_j$, which is the summation of the hidden embedding e_j and the learned attention weight α_j . The final output of the network is predicted by a **value network**, which is defined as

$$\tilde{y}_l := f(a, I_i | \theta) = w^y \phi_v(e_l^a), \quad (8)$$

where θ is the parameters in the defined sequential neural networks. Specifically, ϕ_f , ϕ_i , ϕ_o are sigmoid functions. ϕ_c and ϕ_e are hyperbolic tangent functions. ϕ_a and ϕ_v are rectified linear functions. The proposed **RSN** and **ASN** share the same neural architecture, but use different parameters. That is $\mathcal{R}(a, I_i | \theta_{\mathcal{R}}) = f(a, I_i | \theta_{\mathcal{R}})$ and $\mathcal{S}(a, I_i | \theta_{\mathcal{S}}) = f(a, I_i | \theta_{\mathcal{S}})$. $\mathcal{R}(\cdot, I_i | \theta_{\mathcal{R}})$ and $\mathcal{S}(\cdot, I_i | \theta_{\mathcal{S}})$ denote two vectors of inference value for all $a \in A(I_i)$.

Continual Improvement

With the proposed framework of double neural CFR, it is easy to initialize the neural networks from an existing strategy profile based on the tabular or neural representation. For information set I_i and action a , define $R'_i(a|I_i)$ as the cumulative regret and $S'(a|I_i)$ as the cumulative numerator of average strategy. We can clone the cumulative regret for all information sets and actions by optimizing

$$\theta_{\mathcal{R}}^* \leftarrow \underset{\theta_{\mathcal{R}}}{\operatorname{argmin}} \sum_{i \in N, I_i \in \mathcal{I}_i, a \in A(I_i)} \left(\mathcal{R}(a, I_i | \theta_{\mathcal{R}}) - R'_i(a|I_i) \right)^2. \quad (9)$$

Similarly, the parameters $\theta_{\mathcal{S}}^*$ for cloning the cumulative numerator of average strategy can be optimized in the same way. Based on the learned $\theta_{\mathcal{R}}^*$ and $\theta_{\mathcal{S}}^*$, we can warm start the double neural networks and continually improve beyond the tabular strategy profile.

Overall Algorithm

Algorithm 1 provides a summary of the proposed double neural counterfactual regret minimization algorithm. In the first iteration, if the system warm starts from tabular based CFR or MCCFR methods, the techniques in section Continual Improvement will be used to clone the cumulative regrets and strategy. If there is no warm start initialization, we can start our algorithm by randomly initializing the parameters in RSN and ASN at iteration $t = 1$. Then sampling

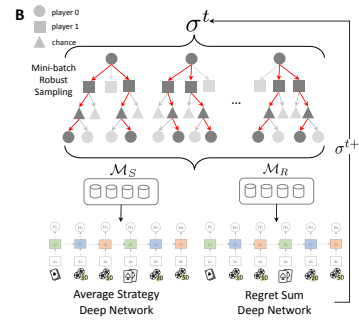


Figure 2: (A) architecture of the sequential neural networks. (B) an overview of the double neural method.

methods will return the counterfactual regret and the numerator of average strategy for the sampled information sets in this iteration, and they will be saved in memories \mathcal{M}_R and \mathcal{M}_S respectively. Then these samples will be used by the NeuralAgent algorithm from Algorithm2 to optimize RSN and ASN. Further details for the sampling methods and the NeuralAgent fitting algorithm will be discussed in the next section.

Algorithm 1: Double Neural CFR

```

1 Function Agent ( $T, b$ ):
2   For  $t = 1$  to  $T$  do
3     if  $t = 1$  and using warm starting then
4       initialize  $\theta_R^t$  and  $\theta_S^t$  from a checkpoint
5        $t \leftarrow t + 1$  ▷ skip cold starting
6     else
7       initialize  $\theta_R^t$  and  $\theta_S^t$  randomly.
8        $\mathcal{M}_R, \mathcal{M}_S \leftarrow$  sampling methods for CFV and
        average strategy.
9       sum aggregate the value in  $\mathcal{M}_R$  by
        information set.
10      remove duplicated records in  $\mathcal{M}_S$ .
11       $\theta_R^t \leftarrow$  Neural( $\mathcal{R}(\cdot|\theta_R^{t-1}), \mathcal{M}_R, \theta_R^{t-1}, \beta_R^*$ )
12       $\theta_S^t \leftarrow$  Neural( $\mathcal{S}(\cdot|\theta_S^{t-1}), \mathcal{M}_S, \theta_S^{t-1}, \beta_S^*$ )
13   return  $\theta_R^t, \theta_S^t$ 

```

Efficient Training

In this section, we will propose three techniques to improve the efficiency of the double neural method. These algorithms can also be used separately in other CFR-family methods.

Robust Sampling Techniques

Theoretically, outcome sampling is more memory efficient than the external sampling, since in outcome sampling only one trajectory is sampled according to strategy profile while in the external sampling, player i will traverse all actions for all information set $I_i \in \mathcal{I}_i$ and the opponent players including chance sample one action. Therefore many information sets will be visited in a sampling process and block $Q_i \in \mathcal{Q}$ will contains many terminal nodes in external sampling. In outcome sampling method, the weighted utility $u_i^{rs}(z) = \frac{u_i(z)}{\pi \sigma_i(z)}$ for terminal node depends on the concrete reach probability $\sigma_i(z)$ in each iteration, therefore it will lead to a high variance and slow down the convergence of the resulting strategy profile.

In this paper, we proposed a new and robust sampling technique which has lower variance than outcome sampling, while being more memory efficient than the external sampling. In this robust sampling method, the sampling profile is defined as $\sigma^{rs(k)} = (\sigma_i^{rs(k)}, \sigma_{-i})$, where player i will randomly select k actions according to sampling strategy $\sigma_i^{rs(k)}(I_i)$ for each information set I_i and other players will randomly select one action according to strategy σ_{-i} .

Specifically, if player i randomly selects $\min(k, |A(I_i)|)$ actions according to discrete uniform distribution $\text{unif}(0, |A(I_i)|)$ at information set I_i , i.e., $\sigma_i^{rs(k)}(a|I_i) = \frac{\min(k, |A(I_i)|)}{|A(I_i)|}$, then

$$\pi_i^{\sigma^{rs(k)}}(I_i) = \prod_{\substack{h \in I_i, h' \sqsubseteq h, \\ h' a \sqsubseteq h, h' \in I_i'}} \frac{\min(k, |A(I_i')|)}{|A(I_i')|} \quad (10)$$

and the weighted utility $u_i^{rs(k)}(z)$ will be a constant number in each iteration, which has a low variance. Because the weighted utility no longer requires explicit knowledge of the opponent's strategy, we can use this sampling method for on-line regret minimization. For simplicity, $k = \max$ refers to $k = \max_{I_i \in \mathcal{I}} |A(I_i)|$ in the following sections.

Lemma 3: If $k = \max$ and $\forall i \in N, \forall I_i \in \mathcal{I}_i, \forall a \in A(I_i), \sigma_i^{rs(k)}(a|I_i) \sim \text{unif}(0, |A(I_i)|)$, then robust sampling is the same as external sampling.

Lemma 4: If $k = 1$ and $\sigma_i^{rs(k)} = \sigma_i$, then robust sampling is the same as outcome sampling.

We provide the proof of **Lemma 3** and **Lemma 4** in (Li et al. 2018).

Mini-batch Techniques

Mini-batch MCCFR: Traditional outcome sampling and external sampling only sample one block in an iteration and provide an unbiased estimation of origin CFV according to **Lemma 2**. In this paper, we present a mini-batch Monte Carlo technique and randomly sample b blocks in one iterations. Let Q^j denote a block of terminals sampled according to the scheme in section Robust Sampling Techniques at j -th time, then **mini-batch CFV** with b mini-batches for information set I_i can be defined as

$$\tilde{v}_i^\sigma(I_i|b) = \sum_{j=1}^b \frac{\tilde{v}_i^\sigma(I_i|Q^j)}{b}. \quad (11)$$

Furthermore, we can show that $\tilde{v}_i^\sigma(I_i|b)$ is an unbiased estimator of the counterfactual value of I_i : **Lemma 5:** $E_{Q^j \sim \text{Robust Sampling}}[\tilde{v}_i^\sigma(I_i|b)] = v_i^\sigma(I_i)$. (we present the proof in (Li et al. 2018).) Similarly, the **cumulative mini-batch regret** of action a is $\tilde{R}_i^T((a|I_i)|b) = \tilde{R}_i^T((a|I_i)|b) + \tilde{v}_i^{\sigma^T}((a|I_i)|b) - \tilde{v}_i^{\sigma^T}(I_i|b)$, where $\tilde{R}_i^0((a|I_i)|b) = 0$. In practice, mini-batch technique can sample b blocks in parallel and help MCCFR to converge faster.

Mini-Batch MCCFR+: When optimizing counterfactual regret, CFR+ (Tammelin 2014) substitutes the regret-matching algorithm (Hart and Mas-Colell 2000) with regret-matching+ and can converge faster than CFR. However, (Burch 2017) showed that MCCFR+ actually converge slower than MCCFR when mini-batch is not used. In our paper, we derive mini-batch version of MCCFR+ which updates cumulative mini-batch regret $\tilde{R}_i^{T,+}((a|I_i)|b)$ up to iteration T by $(\tilde{R}_i^{T-1,+}((a|I_i)|b) + \tilde{v}_i^{\sigma^T}((a|I_i)|b) - \tilde{v}_i^{\sigma^T}(I_i|b))^+$ for $T > 0$. If $T = 0$, $\tilde{R}_i^{T,+}((a|I_i)|b) = (\tilde{v}_i^{\sigma^T}((a|I_i)|b) - \tilde{v}_i^{\sigma^T}(I_i|b))^+$. where $(x)^+ = \max(x, 0)$. In practice, we find that mini-batch MCCFR+ converges faster than mini-batch

MCCFR when specifying a suitable mini-batch size.

Neural Agent for Optimization

Algorithm 2: Optimization of Deep Neural Network

```

1 Function Neural( $f(\cdot|\cdot)$ ,  $\mathcal{M}$ ,  $\theta^{T-1}$ ,  $\beta^*$ ):
2   initialize optimizer, scheduler
3    $\theta^T \leftarrow \theta^{T-1}$ ,  $l_{best} \leftarrow \infty$ ,  $t_{best} \leftarrow 0$ 
4   For  $t = 1$  to  $\beta_{epoch}$  do
5      $loss \leftarrow []$  ▷ initialize loss as an empty list
6     For each training epoch do
7        $\{x^{(i)}, y^{(i)}\}_{i=1}^m \sim \mathcal{M}$  ▷ sampling a mini-batch
8        $batch\_loss \leftarrow \frac{1}{m} \sum_{i=1}^m (f(x^{(i)}|\theta^{T-1}) +$ 
9          $y^{(i)} - f(x^{(i)}|\theta^T))^2$ 
10      back propagation  $batch\_loss$  with
11        learning rate  $lr$ 
12      clip gradient of  $\theta^T$  to  $[-\epsilon, \epsilon]^d$ 
13       $optimizer(batch\_loss)$ 
14       $loss.append(batch\_loss)$ 
15     $lr \leftarrow scheduler(lr)$  ▷ reduce learning rate.
16    if  $avg(loss) < \beta_{loss}$  then
17       $\theta_{best}^T \leftarrow \theta^T$ , early stopping.
18    else if  $avg(loss) < l_{best}$  then
19       $l_{best} = avg(loss)$ ,  $t_{best} \leftarrow t$ ,  $\theta_{best}^T \leftarrow \theta^T$ 
20    if  $t - t_{best} > \beta_{re}$  then
21       $lr \leftarrow \beta_{lr}$  ▷ reset learning rate.
22  return  $\theta^T$ 

```

Define β_{epoch} as training epoch, β_{lr} as learning rate, β_{loss} as the criteria for early stopping, β_{re} as the upper bound for the number of iterations from getting the minimal loss last time, θ^{t-1} as the parameter to optimize, $f(\cdot|\theta^{t-1})$ as the neural network, \mathcal{M} as the training sample consisting information set and the corresponding target. To simplify notations, we use β^* to denote the set of hyperparameters in the proposed deep neural networks. $\beta_{\mathcal{R}}^*$ and $\beta_{\mathcal{S}}^*$ refer to the sets of hyperparameters in RSN and ASN respectively. Algorithm 2 presents the details of how to optimize the proposed neural networks.

Both $\mathcal{R}(a, I_i|\theta_{\mathcal{R}}^{t+1})$ and $\mathcal{S}(a, I_i|\theta_{\mathcal{S}}^t)$ are optimized by mini-batch stochastic gradient descent method. In this paper, we use Adam optimizer (Kingma and Ba 2014) with both momentum and adaptive learning rate. Some other optimizers in (Ruder 2017), however, do not achieve better experimental results. In practice, existing optimizers could not return a relatively low enough loss because of potential saddle point or local minima. To obtain a relatively higher accuracy and lower optimization loss, we use a carefully designed scheduler to reduce the learning rate when the loss has stopped decrease. Specifically, the scheduler reads a metrics quantity, *e.g.* mean squared error, and if no improvement is seen for a number of epochs, the learning rate is reduced by a factor. In addition, we will reset the learning rate in both optimizer and scheduler once loss stops decrease in β_{re} epochs. Gradient clipping mechanism is used to limit the

magnitude of the parameter gradient and make optimizer behave better in the vicinity of steep cliffs. After each epoch, the best parameter will be updated. Early stopping mechanism is used once the lowest loss is less than the specified criteria β_{loss} .

Experiment

The proposed double neural CFR algorithm will be evaluated in No-Limit Leduc Hold'em with stack size 5 and One-Card-Poker game with 5 cards. We will compare it with tabular CFR and deep reinforcement learning based method such as NFSP. The experiments show that the proposed double neural algorithm can converge to comparable results produced by its tabular counterpart while performing much better than deep reinforcement learning method. The current results open up the possibility for a purely neural approach to directly solve large IIG.

Settings. To simplify the expression, the abbreviations of different methods are defined as follows. **XFP** refers to the full-width extensive-form fictitious play method in (Heinrich, Lanctot, and Silver 2015), **NFSP** refers to the reinforcement learning based fictitious self-play method in (Heinrich and Silver 2016). **RS-MCCFR** refers to the proposed robust sampling MCCFR. This method with regret matching+ acceleration technique is denoted by **RS-MCCFR+**. To evaluate the contribution of each neural agent, we replace the tabular based cumulative regret and numerator with RSN and ASN separately. These methods only containing one neural network are denoted by **RS-MCCFR+-RSN** and **RS-MCCFR+-ASN** respectively. **RS-MCCFR+-RSN-ASN** refers to the proposed double neural MCCFR. More specifically, we investigated the following questions.

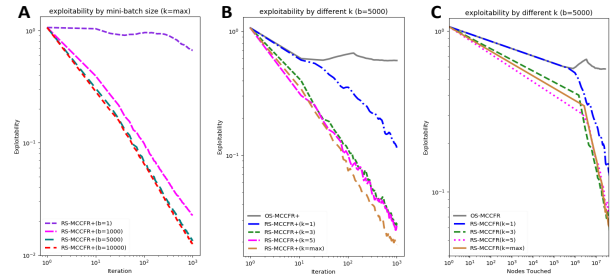


Figure 3: Comparison of different CFR-family methods in Leduc Hold'em. (A) Performance of robust sampling with different batch size. (B) Performance of robust sampling with different parameter k by iteration. (C) Performance by the number of touched node.

Is mini-batch sampling helpful? Figure 3(A) presents the convergence curves of the proposed robust sampling method with $k = max$ under different mini-batch sizes ($b=1, 1000, 5000, 10000$ respectively). The experimental results show that larger batch sizes generally lead to better strategy profiles. Furthermore, the convergence for $b = 5000$ is as good as $b = 10000$. Thus in the later experiments, we set the mini-batch size equal to 5000.

Is robust sampling helpful? Figure 3 (B) and (C) presents convergence curves for outcome sampling, exter-

nal sampling ($k = max$) and the proposed robust sampling method under the different number of sampled actions. The outcome sampling cannot converge to a low exploitability smaller than 0.1 after 1000 iterations (touch more than 10^7 nodes as shown in Figure 3(C) because of the high variance. The proposed robust sampling algorithm with $k = 1$, which only samples one trajectory like the outcome sampling, can achieve a better strategy profile after the same number of iterations. With an increasing k , the robust sampling method achieves an even better convergence rate. Experiment results show $k = 3$ and 5 have a similar trend with $k = max$, which demonstrates that the proposed robust sampling achieves similar strategy profile but requires less memory than the external sampling. We choose $k = 3$ for the later experiments in Leduc Hold'em Poker. Figure 3 (C) presents the results in a different way and displays the relation between exploitability and the cumulative number of touched nodes. The robust sampling with small k is just as good as the external sampling while being more memory efficient on the condition that each algorithm touches the same number of nodes.

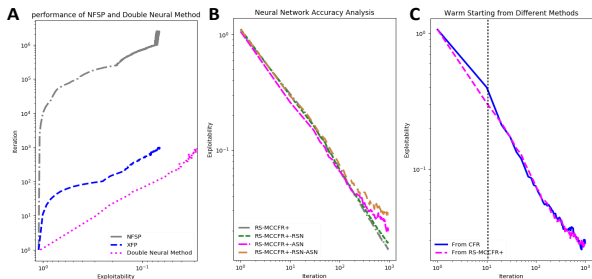


Figure 4: Performance of different methods in Leduc Hold'em. (A) comparison of NSFP, XFP and the proposed double neural method. (B) each contribution of RSN and ASN. (C) continue improvement from tabular methods

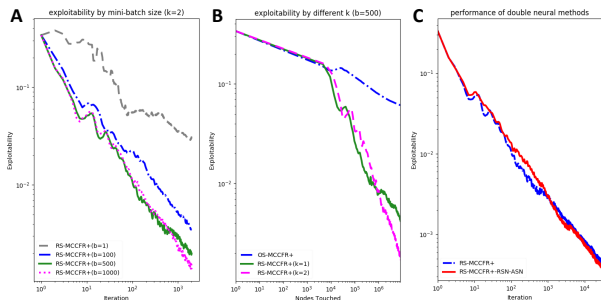


Figure 5: Comparison of different CFR-family methods and neural network methods in One-Card-Poker. (A) Comparison of the robust sampling with different mini-batch size. (B) Comparison of the outcome sampling and the robust sampling with different sample actions k . (C) Comparison of tabular based RS-MCCFR+ and the double neural method.

How does double neural CRF compared to tabular counterpart, XFP and NFSP? To obtain an approximation of Nash equilibrium, Figure 4(A) demonstrates that NFSP

needs 10^6 iterations to reach a 0.06-Nash equilibrium, and requires 2×10^5 state-action pair samples and 2×10^6 samples for supervised learning respectively. The XFP needs 10^3 iterations to obtain the same exploitability, however, this method is the precursor of NFSP and updated by a tabular based full-width fictitious play. Our proposed neural method only needs 200 iterations to achieve the same performance which shows that the proposed double neural algorithm converges significantly better than the reinforcement learning counterpart. In practice, our double neural method can achieve an exploitability of 0.02 after 1000 iterations, which is similar to the tabular method.

What is the individual effect of RSN and ASN? Figure 4(B) presents ablation study of the effects of RSN and ASN network respectively. Both MCCFR+-RSN and MCCFR+-ASN, which only employ one neural network, perform only slightly better than the double neural method. All the proposed neural methods can match the performance of the tabular based method. For RSN, we set the hyperparameters as follows: neural batch size is 256, hidden size is 128 and learning rate $\beta_{lr} = 0.001$. A scheduler, who will reduce the learning rate based on the number of epochs and the convergence rate of loss, help the neural agent to obtain a high accuracy. The learning rate will be reduced by 0.5 when loss has stopped improving after 10 epochs. The lower bound on the learning rate of all parameters in this scheduler is 10^{-6} . To avoid the algorithm converging to potential local minima or saddle point, we will reset the learning rate to 0.001 and help the optimizer to learn a better performance. θ_{best}^T is the best parameters to achieve the lowest loss after T epochs. If average loss for epoch t is less than the specified criteria $\beta_{loss}=10^{-4}$, we will early stop the optimizer. We set $\beta_{epoch} = 2000$ and update the optimizer 2000 maximum epochs. For ASN, we set the hidden size as 256, the loss of early stopping criteria as 10^{-5} . The learning rate will be reduced by 0.7 when loss has stopped improving after 15 epochs. Other hyperparameters in ASN are similar to RSN.

How well does continual improvement work? In practice, we usually want to continually improve our strategy profile from an existing checkpoint (Brown and Sandholm 2016). In the framework of the proposed neural counterfactual regret minimization algorithm, warm starting is easy and friendly. Firstly, we employ two neural networks to clone the existing tabular based cumulative regret and the numerator of average strategy by optimizing Eq. (9). Then the double neural methods can continually improve the tabular based methods. As shown in Figure 4(C), warm start from either full-width based or sampling based CFR the existing can lead to continual improvements. Specifically, the first 10 iterations are learned by tabular based CFR and RS-MCCFR+. The remaining iterations are continually improved by the double neural method.

The proposed double neural method also presents comparable results with tabular-based methods in the experiment of One-Card-Poker as shown in Figure 5.

References

- [Brown and Sandholm 2016] Brown, N., and Sandholm, T. 2016. Strategy-based warm starting for regret minimization in games. 432–438. AAAI.
- [Brown and Sandholm 2017] Brown, N., and Sandholm, T. 2017. Superhuman ai for heads-up no-limit poker: Libratus beats top professionals. *Science* eaa01733.
- [Burch et al. 2012] Burch, N.; Lanctot, M.; Szafron, D.; and Gibson, R. G. 2012. Efficient monte carlo counterfactual regret minimization in games with many player actions. In *Advances in Neural Information Processing Systems*, 1880–1888.
- [Burch 2017] Burch, N. 2017. Time and space: Why imperfect information games are hard. PhD thesis.
- [Cho, Courville, and Bengio 2015] Cho, K.; Courville, A.; and Bengio, Y. 2015. Describing multimedia content using attention-based encoderdecoder networks. arXiv preprint arXiv:1507.01053.
- [Desimone and Duncan 1995] Desimone, R., and Duncan, J. 1995. Neural mechanisms of selective visual attention. Number 18, 193–222. Annual review of neuroscience.
- [Hart and Mas-Colell 2000] Hart, S., and Mas-Colell, A. 2000. A simple adaptive procedure leading to correlated equilibrium. *Econometrica* (65(5)):1127–1150.
- [Heinrich and Silver 2016] Heinrich, J., and Silver, D. 2016. Deep reinforcement learning from self-play in imperfect-information games. arXiv preprint arXiv:1603.01121.
- [Heinrich, Lanctot, and Silver 2015] Heinrich, J.; Lanctot, M.; and Silver, D. 2015. Fictitious self-play in extensive-form games. 805–813. International Conference on Machine Learning.
- [Johanson et al. 2012] Johanson, M.; Bard, N.; Lanctot, M.; Gibson, R.; and Bowling, M. 2012. Efficient nash equilibrium approximation through monte carlo counterfactual regret minimization. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, 837–846. International Foundation for Autonomous Agents and Multiagent Systems.
- [Kingma and Ba 2014] Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- [Lanctot et al. 2009] Lanctot, M.; Kevin, W.; Martin, Z.; and Bowling, M. 2009. Monte carlo sampling for regret minimization in extensive games. In *Advances in neural information processing systems*.
- [Li et al. 2018] Li, H.; Hu, K.; Ge, Z.; Jiang, T.; Qi, Y.; and Song, L. 2018. Double neural counterfactual regret minimization.
- [Moravk et al. 2017] Moravk, M.; Martin, S.; Neil, B.; Viliam, L.; Morrill, D.; Bard, N.; Davis, T.; Waugh, K.; Johanson, M.; and Bowling, M. 2017. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science* (6337):508–513.
- [Osborne and Ariel 1994] Osborne, M. J., and Ariel, R. 1994. *A course in game theory*, volume 1. MIT Press.
- [Ruder 2017] Ruder, S. 2017. An overview of gradient descent optimization algorithms. arXiv preprint arXiv:1609.04747.
- [Tammelin 2014] Tammelin, O. 2014. Solving large imperfect information games using cfr+. arXiv preprint.
- [Waugh et al. 2015] Waugh, K.; Morrill, D.; Bagnell, J. A.; and Bowling, M. 2015. Solving games with functional regret estimation. In *AAAI*, volume 15, 2138–2144.
- [Zinkevich et al. 2007] Zinkevich, M.; Michael, J.; Michael, B.; and Piccione, C. 2007. Regret minimization in games with incomplete information. *Advances in neural information processing systems*.