

Kierunek: **TIN**
Specjalność: **TIP**

PRACA DYPLOMOWA
INŻYNIERSKA

**Implementacja modelu sztucznej inteligencji do gry w
pokera**

Michał Żarejko

Opiekun pracy
Dr. Inż. Paweł Zyblewski

Słowa kluczowe: Deep CFR, Poker, Uczenie przez wzmocnienie, Drzewa decyzyjne

Abstract

The main goal of the thesis is to present author's implementation of the Deep CFR algorithm with the Heads Up Limit Poker Texas Hold'em. It is a modern method for creating artificial intelligence in large partial-observable games. Such environments have always been a great challenge and the main barrier to the development of machine learning. The thesis will show this problem by implementation of Deep CFR and analysis of the results. For this purpose, five recognition models were created every 10 iterations of the algorithm. The next step was to create games with that models. The results allowed to select the best playing models and to see how Deep CFR learned over time. Additionally, the quality of the models was tested against a simple program which simulate beginner player.

Contents

1	Introduction	7
1.1	Genesis of the topic	8
1.2	Goal	8
2	Analiza istniejących rozwiązań	9
2.1	Analiza Texas Hold'em Poker	9
2.2	Uczenie przez wzmacnianie	11
2.3	Teoria Gier	12
2.4	Historia modeli Texas Hold'em Poker	14
2.5	Counterfactual Regret Minimization	15
2.5.1	Regret Matching	15
2.5.2	Counterfactual Regret	15
2.6	Monte Carlo Conterfactual Regret Minimization	18
2.7	Deep CFR	19
2.8	Podsumowanie	21
3	Implementacja algorytmu	23
3.1	Implementacja sieci neuronowych	24
3.1.1	Architektura modelu	25
3.1.2	Budowa zbiorów danych	26
3.1.3	Proces uczenia	27
3.2	Implementacja środowiska	28
3.3	Implementacja Deep CFR	29
3.4	Podsumowanie	30
4	Wyniki	31
4.1	Proces uczenia modeli rozpoznawania	31
4.2	Wyniki rozgrywek modeli	35
4.3	Porównanie modeli z programem HP	38
4.4	Podsumowanie	41
5	Wnioski i podsumowanie etapów pracy	42
5.1	Wnioski	43
5.2	Dalszy rozwój algorytmów bazujących na metodzie CFR	44

Chapter 1

Introduction

Machine learning has been in development for a long time. As early as the 40s, books and programs related to artificial intelligence [1] have been writting. In 1940 Donald Hebb created theoretical foundations used in modern neural networks [1]. Despite such early progress, the branch of science, has begun to achieve great success about 10 years ago. This is due to the fact that many algorithms need a lot of computational power [1].

Today, there are many tools under development which are related with this topic. For example voice assistants, language translators, models to display elements on websites, video games or smart cars. Additionally, artificial intelligence is heavily used in companies, on production halls, in transportation, medicine or cyber security.

Despite many possibilities and applications, AI is currently gaining the most popularity media through competitive games, where the main task is to show the superiority of algorithms over humans. Today, there are many such events, where world champions in a given game, lost to recognition models.

In 2016, there was an organized match between Fan Hui, the European champion in the Chinese game Go and the algorithm AlphaGo. The model created by the DeepMind team won against the opponent. Until now, the game was widely considered to be complicated and difficult to solve.

In 2019, a model OpenAI Five as first AI in the world, beat Team OG in the e-sport competition in the game Dota 2. The event was heavily covered in the media because of the first such achievement in this sport. Additionally, Dota 2 was a very complex environment. For example, the game Go solved a few years earlier, contained 150 possible moves per turn, Dota 2 could have 20,000 of them in less than an hour [2].

There are many such events. They show that today's artificial intelligence can surpass humans in strategic thinking. Additionally, it prove that AI topic is still in developed with more and more interest each year.

1.1 Genesis of the topic

Often in developing ML programs, a major challenge is the level of complexity of the game. This depends on whether the environment is deterministic or stochastic, how dynamic game is or whether the dimensional space is discrete or infinite. Currently, however, one of the major problems of such programs is the insufficient extent of available information about the environment [5].

Artificial intelligence, in order to win, needs a large number of unambiguous inputs for subsequent learning. An example of a game that meets such requirements is chess. AI makes moves based on information such as the arrangement of pawns in a given turn along with previous settings. The change in the environment on the opponent's side is immediately visible so model can easier bind actions and observations.

An example of a game hard to teach AI that does not meet this condition is Poker Texas Hold'em. It exhibits an incomplete set of information, despite knowledge of the cards in the hand and on the table, the player has no knowledge of opponent cards. In this case, two seemingly identical states of the environment in reality may differ. In addition, the environment is characterized by very high randomness which is hard to predict good moves. Because of these characteristics, most popular algorithms such as DQN (*Deep Q Learning*), Actor-Critic or AlphaZero become useless and do not give good results.

This paper presents a way to possibly solve such a problem using the Deep CFR algorithm [3]. This is a popular method for creating models of recognition in card games.

1.2 Goal

The main goal of this paper is to implement a Deep CFR algorithm that will create 5 models of recognition in the game HULH (*Heads Up Limit Texas Poker Hold'em*). This is a popular version of 2-player play where participants cannot choose the amount of the raise on their own. It is limited by a set value. This environment minimizes possible moves to 3 actions, making it a simpler base for machine learning. The resulting models will then be used to create games consisting of all combinations of the two models, where each game will be repeated 200 times. The second stage will be to calculate the average pool won and lost by each AI along with the distribution of moves made. This process will determine which model performs best and what strategy it uses to play the game.

The work has been divided into 5 stages for this purpose. The next section is a theoretical introduction to the implementation of the program. It describes possible solutions to the problem, analysis of the game Poker Texas Hold'em and the required theory for understanding Deep CFR. The third section presents the written program along with the technologies used. The final two chapters discuss the learning results, the results model games, and a summary of the work.

Chapter 2

Analysis of existing solutions

W ciągu ostatnich 15 lat powstało wiele algorytmów rozwiązujących różne wersje gry Poker. Miedzy innymi CFR (*Counterfactual Regret Minimization*) [5], XFP (*Extensive-Form Fictitious Play*) [4] lub NFSP (*Neural Fictitious Self-Play*) [6]. Pierwszy z wymienionych, CFR powstał w 2007 roku. Był pomyslną próbą rozwiązania abstrakcyjnego środowiska Poker Texas Hold'em [5]. Na jego podstawie utworzono wiele nowoczesnych algorytmów, które dają szansę rozwiązać takie gry jak HULH [5].

Z wymienionych metod zaimplementowanym rozwiązaniem w niniejszej pracy jest CFR rozszerzony o sieci neuronowe, czyli Deep CFR z gra HULH. Pozwala on na szybsze trenowanie modeli w środowisku typu zero-sum, dodatkowo lepiej rozwiązuje gry o dużych rozmiarach [3].

W tym rozdziale zostana przedstawione profesjonalne sposoby wyboru strategii w grze Poker Texas Hold'em. Określa one cechy, jakimi powinien charakteryzować się prawidłowo utworzony model rozpoznawania. Następnie rozdział przedstawi istniejące sztuczne inteligencje, które wykorzystały podstawy metody CFR do zwyciężania z profesjonalnymi graczami. Ostatnim etapem jest teoria związana z algorytmem Deep CFR.

2.1 Analiza Texas Hold'em Poker

Jest to jedna z najpopularniejszych gier rywalizacyjnych w kasynach, dodatkowo jest to dominująca gra hazardowa. Można ją scharakteryzować dużą stochastycznością oraz częściową obserwowalnością, tab. 2.1 [7].

Przez takie cechy gra była od zawsze tematem sporów, czy na jej wynik ma większy wpływ losowość, czy umiejętność. Dużym aspektem pomagającym w osiągnięciu zwycięstwa jest panowanie nad emocjami w celu ukrycia informacji o posiadanych kartach przez przeciwnikiem. Druga zasada jest umiejętność decydowania kiedy grać agresywnie, a kiedy pasywnie.

Table 2.1: Charakterystyka wybranych gier

	środowisko deterministyczne	środowisko niedeterministyczne
pełny zestaw informacji	szachy Go	Monopoly Tetris
niepełny zestaw informacji	Saper Mahjong	Poker Makao

Dodatkowo ważnym elementem jest obserwacja gry oraz wybór prawidłowej strategii. Można ją wybrać na bazie dostępnych kart i dotychczasowego zachowania oponenta.

Kolejna ważna zasada jest obserwacja przeciwnika oraz zapamiętywanie jego poprzednich akcji. Przez to można określić, czy gra on agresywnie, czy pasywnie i dobrać do niego odpowiednią strategię. W tym celu dokonuje się klasyfikacji przeciwników na bazie częstotliwości wykonywanych ruchów [8].

Każdego gracza można podzielić na cztery grupy, Loose Aggressive, Loose Passive, Tight Passive oraz Tight Aggressive [8]. Prawidłowe rozpoznanie danego stylu gry może zdecydować o wyborze prawidłowej strategii i zwycięstwie. Poniżej opisano każdy z nich oraz porównano je z rys. 2.2, który pokazuje, w jakim stopniu profesjonalny gracz powinien używać każdego z nich [8].

Loose Passive

Osoba, która bardzo często wchodzi do gry niezależnie czy posiadane karty dają jej wysokie szanse na wygraną. Taki sposób grania nie jest dobrą strategią, ponieważ można łatwo się do niego dostosować przez używanie tylko mocnych kart [8]. Taki styl jest często używany przez niedoświadczonych graczy.

Loose Aggressive

Osoba, która często przebija stawke nawet w rundzie *pre-flop* i często wchodzi do gry [8]. Strategia ma stworzyć przekonanie wśród oponentów, że gracz ma bardzo duże szanse wygranej od samego początku. Okazuje się ona jednak nieefektywna, jeśli przeciwnicy nie pasują w początkowych etapach [8].

Tight Passive

Uczestnik gry wchodzi tylko z dobrymi kartami, wykonując często akcje *call*. Ostatecznie pasuje przy spotkaniu z graczem agresywnym. Taka osoba gra bardzo dokładnie tak, aby mało ryzykować, przez co często traci wiele okazji gdzie mogła by wygrać [8]. W rundach

traci ona bardzo małe stawki ale pomimo tego, bardzo rzadko wygrywa. W rezultacie stopniowo traci żetony przy grze z bardziej doświadczonymi graczami.

Tight Aggressive

Graja podobnie do typu Tight Passive w początkowych etapach gry, a następnie zmieniają swój styl na mniej lub bardziej agresywny tak aby wygrać rundę po mimo większego ryzyka [8]. Rys. 2.1 pokazuje, że jest to najlepsza wersja strategii, jaka można grać, dlatego często jest ona wybierana przez profesjonalnych graczy.

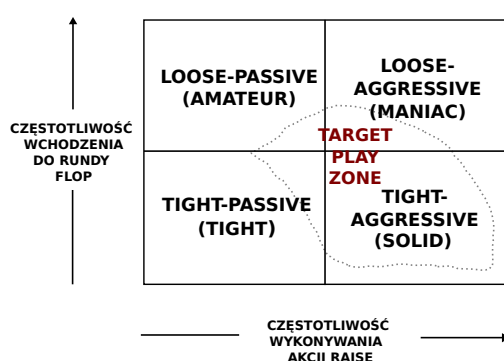


Figure 2.1: Podział graczy w grze Poker Texas Hold'em [8].

Jak wynika z powyższych kategorii, gra Poker Texas Hold'em zawiera wiele elementów niezwiązanych z losowością, gdzie obserwacje i dobieranie odpowiedniej strategii do typu gracza pełni kluczową funkcję. Korzystając z tych informacji, będzie można określić poziom zaawansowania i styl gry utworzonych modeli, które są opisane w rozdziale 4.

2.2 Uczenie przez wzmocnianie

Jest wiele sposobów na stworzenie sztucznych inteligencji w grach. Między innymi można użyć technik uczenia nadzorowanego pod warunkiem, jeśli przygotuje się odpowiednie zbiory danych. W pracy jednak zdecydowano się na stworzenie modelu, korzystając z uczenia przez wzmocnianie, czyli rozwiązania gdzie AI nie potrzebuje wstępnej bazy uczącej. Wynika to z faktu, że jest mało publicznych zapisów gry Poker Texas Hold'em z profesjonalnymi graczami, które mogłyby posłużyć jako zbiory danych.

Algorytmy należące do wybranego działu powinny być w stanie polepszać swoje wyniki na podstawie interakcji ze środowiskiem bez korzystania z zewnętrznych materiałów. Wiele istniejących metod należących do wybranej techniki zakłada, że środowisko można opisać przez model matematyczny MDP (*Markov Decision Process*). Określa ona sekwencyjnie podejmowane decyzje wraz z uzyskiwaną nagrodą od środowiska [9]. W każdym z nowych

stanów, w jakich znajduje się agent (AI), wykonuje on pojedyncze akcje a , zyskując od środowiska informacje o nowym stanie s oraz nagrodzie r . Elementem wyjściowym zasady powinien być zbiór strategii w postaci modelu, rys. 2.2.

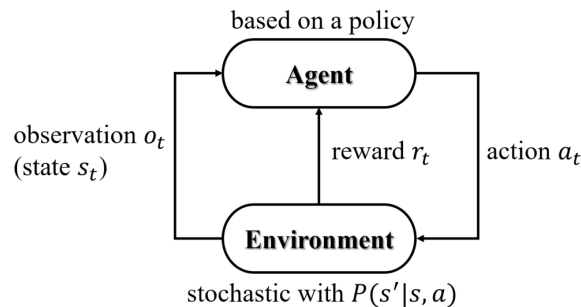


Figure 2.2: Schemat modelu matematycznego MDP [9].

MDP może opisywać jedynie środowiska z pełnym zakresem informacji, w przypadku algorytmu będącego tematem pracy, głównym zadaniem jest rozwiązanie środowiska z niepełnym zestawem danych. Wtedy należy rozpatrywać zasadę POMDP (*Partially Observable Markov Decision Process*) [9].

W przeciwieństwie do poprzedniego modelu w tym przypadku agent nie zna aktualnego stanu, w którym się znajduje [9]. Przez takie okoliczności musi połączyć zależności wykonywane akcje i obserwacje, a nie stany. Większość gier karcianych można zakwalifikować do tego typu problemów [9].

2.3 Teoria Gier

Aby zrozumieć działanie algorytmów CFR, Deep CFR, MCCFR itd. należy zapoznać się z podstawami działu matematyki o nazwie Teoria Gier. Bada on optymalne zachowanie w środowiskach, gdzie występuje konflikt [10]. W przypadku gry Poker Texas Hold'em zostały opisane terminy Równowaga Nasha oraz postać ekstensywna. Są one obowiązkowe do zrozumienia algorytmu Deep CFR.

Gra w postaci ekstensywnej

Gry w formie ekstensywnej można przedstawić jako drzewo decyzyjne, gdzie każdy węzeł rozgałęzia się na możliwe akcje oraz identyfikuje aktualny stan gracza przez zestaw informacji. Stany końcowe drzewa określają zysk lub stratę nagrody wybranego gracza [3]. Jest to sposób na uproszczenie opisu gry, gdzie ruchy są wykonywane nierównocześnie.

Na rys. 2.4 przedstawiono przykład gry 'Papier-Kamień-Nożyce' w formie ekstensywnej, gdzie gracze P1 i P2 eksplorują 3 akcje w swoich węzłach. Każdy z nich uzyskuje wyniki danych ścieżek oraz zapamiętuje dotychczasową historię, co może zostać potem wykorzystane do znalezienia najbardziej opłacalnych strategii.

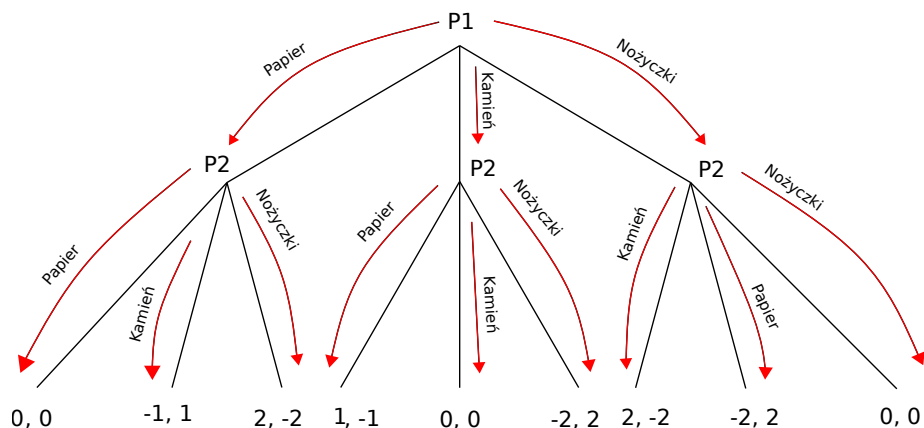


Figure 2.3: Opis gry w postaci ekstensywnej.

Równowaga Nasha

W grach to twierdzenie określa perfekcyjny stan wyboru akcji, gdzie wszyscy gracze wykorzystują najlepszy zestaw strategii, którego zmiana przyniesie tylko straty. Oznacza to, że nie jest możliwa zmiana ruchów na lepsze oraz zwiększenie uzyskanej nagrody przy spełnieniu tej zasady [10].

Dobrym przykładem prezentującym taki stan jest "Dylemat Więźnia" [10]. Takie środowisko zawiera dwóch przestępców, którzy są przesłuchiwanymi w odseparowanych pomieszczeniach co oznacza, że każdy z nich ma ograniczony zestaw informacji. Mają oni dwie opcje, przyznać się do zarzutów lub tego nie robić. Każda z kombinacji akcji jest zaprezentowana w tab. 2.2, gdzie wartości określają lata spędzone w więzieniu po danym ruchu. Posługując się Równowagą Nasha, można stwierdzić, że najlepsza opcja dla obu uczestników będzie przyznawanie się za każdym razem [10]. Wynika to z faktu, że uczestnicy wtedy nie ryzykują.

Głównym zadaniem większości algorytmów gier karcianych bazujących na metodzie CFR jest znalezienie takiego stanu. Deep CFR odkrywa zbiory strategii, które są bliskie Równowadze Nasha w grach karcianych [3].

Table 2.2: Wyniki akcji środowiska "Dylemat więźnia".

	wieźnia A się przyznaje	wieźnia A kłamie
wieźnia B się przyznaje	1	5
wieźnia B kłamie	5	0

2.4 Historia modeli Texas Hold'em Poker

Bazując na Teorii Gier oraz innych twierdzeniach powstało wiele rozwiązań różnych wersji gry Poker. Pierwsze dokumenty naukowe omawiały bardzo proste środowiska rozwiązywane przez algorytmy jak CFR [3].

Dopiero w 2015 roku utworzono znana sztuczna inteligencja "Cepheus" rozwiązująca problem HULH [11]. Było to pierwsze takie osiągnięcie w historii. Kolejnym etapem były prace nad algorytmami mogącym rozwiązać problem gry HUNH (*Heads Up No-limit Texas Hold'em*). Pierwszy powstały model w 2017 roku nazwano "DeepStack" [13]. Mieszał on sieci neuronowe z technikami algorytmu CFR. W podobnym czasie utworzono kolejne, bardziej zaawansowane AI, Libratus [12]. W 2019 roku pierwszy raz udało się rozwiązać problem gry HUNH składającej się z 6 graczy. Zbudowano model o nazwie "Pluribus", który był pierwszym AI rozwiązującym standardową wersję pokera [19].

Pomimo takich osiągnięć tworzone modele do 2019 potrafiły grać tylko w środowiskach składających się maksymalnie z 2 osób typu zero-sum [12] [11] [13]. Wynika to z poziomu skomplikowania gier częściowo-obserwowalnych. Sposoby na jego rozwiązanie zaczęły powstawać od niedawna, a pierwsze dwa duże osiągnięcia w grze HUNH miały miejsce dopiero w 2017 roku.

Cepheus

AI powstałe w celu wygrywania w grach HULH. Był to pierwszy model rozwiązujący dużą wersję gry Texas Poker Hold'em w historii [11]. Wykorzystał on nowszą wersję techniki CFR, która nazwano CFR+ [11]. W wyniku dwóch miesięcy nauki i testów nowa metoda zbiegała się znacznie szybciej do Równowagi Nasha niż bazowy CFR [11]. Powstały model jest udostępniony publicznie, każdy może go przetestować.

DeepStack

Model DeepStack rozwiązał HUNH przez połączenie metody CFR, sieci neuronowych wraz z dodatkowymi elementami. W rezultacie AI zaczęło osiągać bardzo dobre wyniki. Przetestowano go na 33 profesjonalnych graczach w wielu iteracjach gry. Model w większości przypadków wygrał [13]. Była to pierwsza wygrana AI z człowiekiem w normalnej wersji gry Poker Texas Hold'em z taką częstotliwością.

Libratus

Sztuczna inteligencja, która jest wykorzystywana w grach HUNH. Jak wynika z testów wygrywa znacznie częściej niż DeepStack z profesjonalnymi graczami [12]. Przetestowano go z mistrzami gry, Dong Kim, Dan McAulay, Jimmy Chou i Jason Les [12]. AI wygrało z nimi, z ogromną przewagą.

2.5 Counterfactual Regret Minimization

2.5.1 Regret Matching

Jest to nieodłączna metoda uczenia AI w grach karcianych. Polega ona na liczeniu najlepszej strategii pod warunkiem, że znany jest wektor żalu (*regrets*) w węźle [5]. Taki wektor opisuje się jako tablice wag o długości równej liczbie możliwych akcji gracza. Każda z tych wag opisuje, jak dużym błędem będzie niewykonania danego ruchu.

Poniżej przedstawiono wzór wynikający z tej metody, gdzie $R^t(a)$ jest omawianym wektorem [5]. Następnie, aby uzyskać nową strategię, usuwa się wartości ujemne, czyli takie, których gracz nie żałował (formuła nr 2.2). Potem sprawdzana się, czy ich suma jest większa od zera w celu wybrania odpowiedniego wzoru, formuła 2.1. W zależności od tego warunku otrzymuje się określony rozkład prawdopodobieństwa wykonania każdej z akcji.

$$p_i^t(a) = \begin{cases} \frac{R^{T,+}(a)}{\sum_{a' \in A} R^{T,+}(a')} & \text{if } \sum_{a' \in A} R^{T,+}(a') > 0; \\ \frac{1}{|A|} & \text{otherwise.} \end{cases} \quad (2.1)$$

$$R^{t,+}(a) = \max(R^t(a), 0) \quad (2.2)$$

Proces ten jest powtarzany wielokrotnie, tak, aby przy każdej iteracji rozkłady prawdopodobieństwa ruchów były stopniowo poprawiane.

W przypadku algorytmu Deep CFR zachodzi modyfikacja formuły nr 2.1. Strategia jest liczona na dodatnich wartościach żalu podzielonego przez prawdopodobieństwo dostania się do tego stanu $D^T(I, a)$ [3]. Jeśli suma jest ujemna, to zostaje wybrana akcja z najwyższą wartością $D^T(I, a)$ [3].

$$\sigma_i^{t+1}(I, a) = \begin{cases} \frac{D^{T,+}(a)}{\sum_{a' \in A} D^{T,+}(a')} & \text{if } \sum_{a' \in A} D^{T,+}(a') > 0; \\ \operatorname{argmax}(D^T(I, a)) & \text{otherwise.} \end{cases} \quad (2.3)$$

2.5.2 Counterfactual Regret

Algorytm CFR do znanych wcześniej metod dodał termin 'Immediate Counterfactual Regret' oznaczany przez $R_{i,imm}^T(I)$, czyli żal przydzielony do węzła I. Do obliczenia takiego parametru została zdefiniowana wartość "counterfactual utility" $u_i(\sigma, I)$. Oznacza ona przewidywa opłacalność stanu I gdzie wszyscy gracze używają strategii σ [5]. Dodatkowo $\pi^\sigma(h, h')$ oznacza prawdopodobieństwo dostania się z historii h do nowego stanu h' przy strategii σ [5].

$$u_i(\sigma, I) = \frac{\sum_{h \in I, h' \in Z} \pi_{-i}^\sigma(h) \pi^\sigma(h, h') u_i(h')}{\pi_{-i}^\sigma(I)} \quad (2.4)$$

Na podstawie równania 2.5 można wyliczyć końcową wartość żalu w algorytmie CFR.

$$R_{i,\text{imm}}^T(I, a) = \frac{1}{T} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) (u_i(\sigma^t|_{I \rightarrow a}, I) - u_i(\sigma^t, I)) \quad (2.5)$$

Powyższe 2 równania można doprowadzić do formuły nr 2.6. Wartość $\pi^\sigma(h, h')$ została zastąpiona przez liczbę 1, ponieważ CFR zakłada, że dla $u_i(\sigma^t|_{I \rightarrow a}, I)$, gracz wykonuje zawsze akcje a z całkowitą pewnością [5].

$$R_{i,\text{imm}}^T(I, a) = \frac{1}{T} \sum_{t=1}^T \pi_{-i}^{\sigma}(h) \sum_{h \in I, h' \in Z} (1 * u_i(h') - \pi^\sigma(h, h') u_i(h')) \quad (2.6)$$

Po uzyskaniu $R_{i,\text{imm}}^T(I, a)$ można wykorzystać metodę "Regret Matchning" i zaktualizować strategię. Poniżej przedstawiono przykład obliczeń pojedynczego węzła P2 oraz wyniki dla gry "Papier-Kamień-Nożyce" przy ustawionych nagrodach i karach w stanach końcowych jak na rys. 2.4, 2.5, 2.6.

$$\pi^\sigma(h, h') u_i(h') = \left(\frac{1}{3} \cdot 0\right) + \left(\frac{1}{3} \cdot -1\right) + \left(\frac{1}{3} \cdot 2\right) = \frac{1}{3}$$

$$T * R_{i,\text{imm}}^T(I, a) = \left(\left(0 - \frac{1}{3}\right), \left(-1 - \frac{1}{3}\right), \left(2 - \frac{1}{3}\right)\right) \cdot \frac{1}{3} = \left(-\frac{1}{9}, -\frac{4}{9}, \frac{5}{9}\right)$$

$$\sigma_i^{t+1}(I, a) = (0, 0, 1)$$

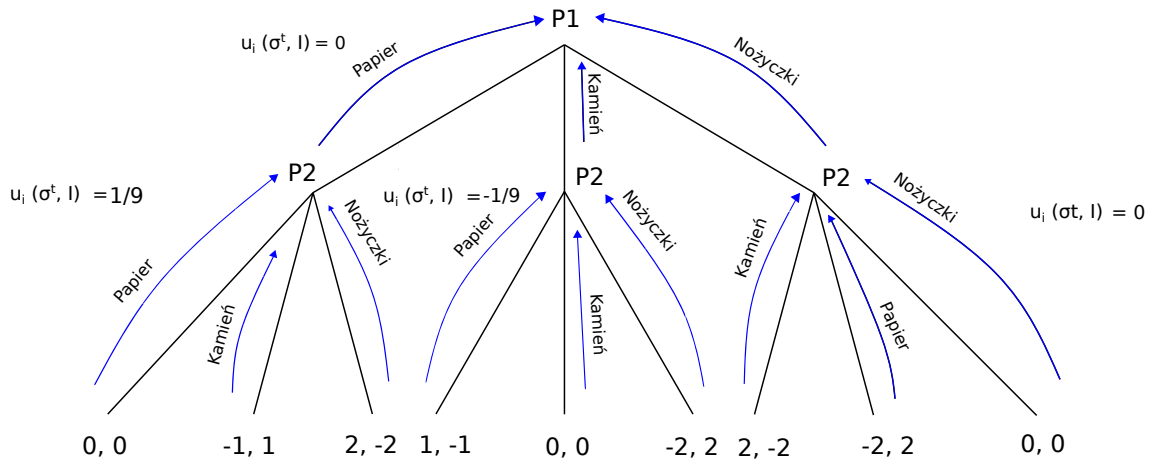


Figure 2.4: Przykład obliczonych wartości 'counterfactual utility'.

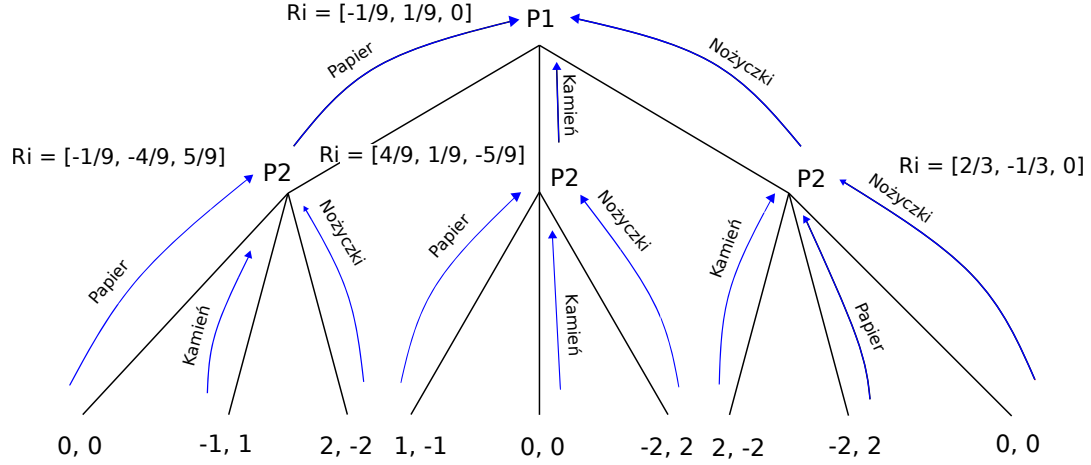


Figure 2.5: Przykład obliczonych wartości 'Immediate Counterfactual Regret'.

Na podstawie rys. 2.4 widać, że gracz P2 będzie miał stan o najwyższej wartości "counterfactual utility" w węźle $u_i(\sigma^t, I) = \frac{1}{9}$, a najniższej dla $u_i(\sigma^t, I) = -\frac{1}{9}$.

Powyżej zaprezentowano wektory $R_{i,\text{imm}}^T(I, a)$. Gracz P1 grając, najbardziej będzie żałował nie wykonania akcji *Kamień*, a najmniej *Papier*.

Rys. 2.6 przedstawia wektory określające jakimi rozkładami akcji powinni się kierować gracze, aby osiągnąć najlepsze wyniki. Są to wektory wyliczone po pierwszej iteracji, w praktyce eksploracja drzewa i powyższe obliczenia są powtarzane wielokrotnie. Kończącym etapem algorytmu CFR jest policzenie średniej strategii, która ma reprezentować Równowagę Nasha [5].

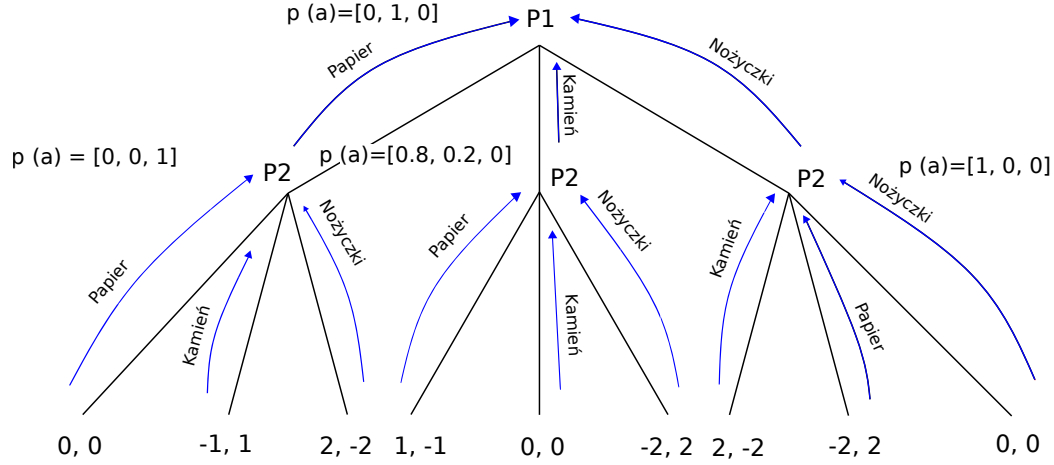


Figure 2.6: Przykład obliczonych strategii.

2.6 Monte Carlo Counterfactual Regret Minimization

Algorytm CFR eksploruje całe drzewa decyzyjne w jednej iteracji, co tworzy wymagania na dużą moc obliczeniową i długi czas uczenia. Dla małych gier takie rozwiązanie jest akceptowalne, ale w przypadku większych środowisk jest nieefektywne. Spowodowało to powstanie nowszej wersji algorytmu CFR, MCCFR (*Monte Carlo Counterfactual Regret Minimization*), który w każdej iteracji eksploruje tylko część drzewa [14].

Metode można podzielić na dwie odmiany, Outcome-Sampling oraz External-Sampling [14].

W pracy został przedstawiony sposób MCCFR ES (*Monte Carlo Counterfactual Regret Minimization External Sampling*), ponieważ taki jest zaimplementowany w algorytmie Deep CFR. MCCFR ES przed eksploracją drzewa wybiera kolejno spośród graczy jednego uczestnika, którego oznacza się jako "traverser" [3]. Eksploruje on wszystkie odpowiedzi ze swoich akcji w danym węźle. W międzyczasie inni uczestnicy wykonują pojedynczy ruch na podstawie swojej najlepszej strategii [14].

Na rys. 2.7 przedstawiono przykład części drzewa HULH, gdzie gracz P2 został wybrany jako "traverser". Jak można zauważyć, tylko jego węzły rozgałęziają się na wszystkie możliwe ścieżki. Dodatkowo w drodze powrotnej obliczono $u_i(\sigma, I)$ dla wszystkich stanów używając wzoru 2.5.

Jest to przykład z jednej eksploracji gry, w praktyce powtarza się ten proces wielokrotnie,

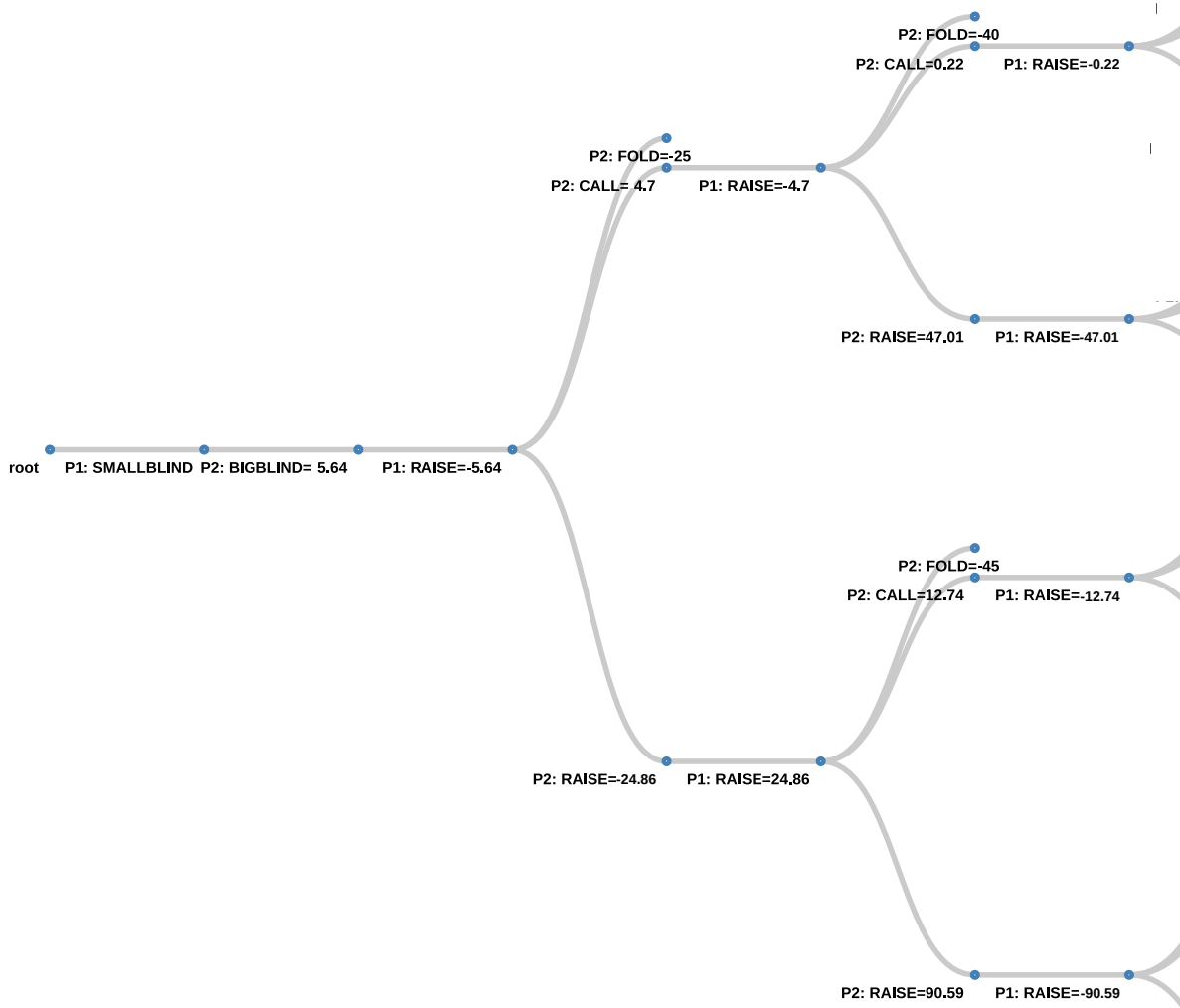


Figure 2.7: Cześć drzewa decyzyjnego MCCFR ES z graczem P2 jako "traverser".

uzyskując różne wersje drzew.

Mając takie wyniki, można przystąpić do liczenia $R_{i,imm}^T(I)$ oraz poprawiania wartości przez "Regret Matching". MCCFR ES spełnia swoją funkcję, ale wymaga wielu iteracji, aby uzyskać dobre wyniki. Z tego powodu w dalszym rozdziale zostanie przedstawiona metoda Deep CFR, która przyspiesza proces uczenia przez użycie sieci neuronowych.

2.7 Deep CFR

Algorytm Deep CFR opisany w [3] rozwija podstawową wersję metody CFR o sieci neuronowe. Taka modyfikacja była wymagana, aby utworzyć AI, które może rozwiązać duże gry jak HULH znacznie szybciej. Liczy on wektory w drzewie decyzyjnym przez algorytm MCCFR ES. Dodatkowo Deep CFR okazał się lepszy niż popularny algorytm NFSP z 2016 roku [3].

Deep CFR wykorzystuje sieci neuronowe do przewidzenia wartości $D^T(I, a)$ w podobnych

obserwacjach, potem na podstawie predykcji liczy strategię ze wzoru nr 2.3. Następnie liczona jest zaktualizowana wersja wektora żalu, formuła nr 2.5. Obliczone wyniki są dodawane do buforów $(B_1, B_2) \in B_p$, a strategia do zbioru B_s .

Po wielu iteracjach rozpoczyna się nauka sieci z zebranej bazy. Poniżej przedstawiono dokładny opis Deep CFR.

Algorithm 1: Deep CFR	
Wejście: $B_p, B_s, \theta_s, \theta_p, P, N, K$	
Wyjście: θ_s	
1 for n in N do	
2 $h, t \leftarrow$ Nowa gra	
3 for p in P do	
4 for k in K do	
5 $\text{MCCFRES}(\theta_p, \theta_{p-1}, p, t, h, B_p, B_s)$	
6 $\theta_p \leftarrow \text{TRAIN}(B_p, \theta_p, t)$	$\triangleright \text{loss: } \frac{1}{N_{\text{batch}}} \sum (y_i - \hat{y}_i)^2 \cdot t_i$
7 $\theta_s \leftarrow \text{TRAIN}(B_s, \theta_s, t)$	$\triangleright \text{loss: } \frac{1}{N_{\text{batch}}} \sum (y_i - \hat{y}_i)^2 \cdot t_i$
8 return θ_s	

Algorytm na wejściu dostaje argumenty przystosowane do gry 2-osobowej. Pierwszymi elementami są bufor graczy (B_1, B_2) oraz kontener na strategię B_s . Dodatkowo metoda potrzebuje listę uczestników P , z których będzie wybierany cyklicznie gracz eksplorujący drzewo decyzyjne. Ostatnimi argumentami są iteracje gry N oraz liczba cykli K metody MCCFR ES.

Podsumowując należy ustawić trzy petle wraz z nową rundą, uzyskując początkową historię oraz krok gry t . Kolejnym etapem jest wykonanie funkcji *MCCFRES*. Przyjmuje ona na wejściu sieci neuronowe obu graczy, listę uczestników p , numer kroku t , historię rundy h oraz bufor. Po k powtórzeniach następuje uczenie sieci neuronowej wybranego wcześniej gracza eksplorującego drzewo.

Model używa zmodyfikowanej wersji funkcji MSE (*Mean Square Error*) do liczenia błędu predykcji. Każdy wynik $(y_i - \hat{y}_i)^2$ mnożony jest przez krok t_i w, którym uzyskano y_i [3]. Dodatkowo jak wynika z badań, algorytm osiąga lepsze wyniki jeśli sieci neuronowe są trenowane od początku [3]. Dlatego należy przed funkcją *TRAIN* wyczyścić model.

Po wszystkich powtórzeniach i zebraniu całej bazy bufora B_s , rozpoczyna się trenowanie sieci neuronowej θ_s w ten sam sposób jak inne modele. Elementem wyjściowym Deep CFR jest sieć θ_s .

Algorithm 2: Implementacja MCCFRES korzystając z sieci neuronowych

```
1 Function MCCFRES( $\theta_p, \theta_{p-1}, p, t, h, B_p, B_s$ ):  
2    $t_r \leftarrow h$  ▷ sprawdzenie czyja jest aktualnie tura  
3   if  $h$  jest stanem końcowym  $Z$  then  
4     return  $u_p(h)$   
5   else if  $p = t_r$  then  
6      $\hat{D}(I) \leftarrow$  utworzenie wektora  $\hat{D}(I)$  z sieci  $\theta_p$   
7      $\sigma(I) \leftarrow$  obliczenie strategii  $\sigma_{t_r}(I)$ , korzystając z  $\hat{D}(I)$  i wzoru 1.1  
8     for  $a$  in  $A(h)$  do  
9        $u_{t_r}(ha) \leftarrow$  MCCFRES( $\theta_p, \theta_{p-1}, p, t+1, h+a, B_p, B_s$ )  
10       $u_{t_r}(\sigma, I) \leftarrow \sum(u_{t_r}(ha) \cdot \sigma(I))$   
11       $R_{t_r}^T(I) \leftarrow (u_{t_r}(ha) - u_{t_r}(\sigma, I))$   
12       $B_{t_r} \leftarrow$  dodanie próbki  $[R_{t_r}^T(I), h, t]$  do bufora  
13   else  
14      $\hat{D}(I) \leftarrow$  utworzenie wektora  $\hat{D}(I)$  z sieci  $\theta_p$   
15      $\sigma(I) \leftarrow$  obliczenie strategii  $\sigma_p(I)$ , korzystając z  $\hat{D}(I)$  i wzoru 1.1  
16      $B_s \leftarrow$  dodanie próbki  $[\sigma(I), h, t]$  do bufora  
17      $a \leftarrow \sigma(I)$   
18     return MCCFRES( $\theta_p, \theta_{p-1}, p, t+1, h+a, B_p, B_s$ )
```

Implementacja MCCFRES przedstawiona powyżej przy zadanych argumentach rozpoczyna się od sprawdzenia, który gracz rozpoczyna daną turę. Na podstawie tej informacji będzie wykonywana dalsza część algorytmu.

Pierwszym krokiem jest sprawdzenie, czy stan gry jest końcowym etapem. W przypadku prawdziwego warunku zwracana jest wygrana lub przegrana wartość stawki. Jeśli powyższy etap jest fałszywy, algorytm sprawdza, czy gracz jest oznaczony jako "traverser". Wtedy program używając sieci neuronowej gracza, otrzymuje wektor $\hat{D}(I)$ przez wprowadzenie do modelu informacji o widocznych kartach oraz dotychczasowej historii gry. Korzystając ze wzoru 2.3, liczy strategię, odczytuje wektor $u_{t_r}(ha)$ wykonując rekurencję. Ostatnim krokiem jest wyliczenie wektora żalu i dodanie go do bufora.

Jeśli powyższy warunek był fałszywy, gracz liczy nowy rozkład akcji i dodaje go do zbioru strategii. Następnie korzystając z tej sieci neuronowej i otrzymanej dystrybucji wykonuje nową akcję.

2.8 Podsumowanie

Środowiska z niepełnym zestawem informacji i brakiem deterministyczności są trudne do rozwiązania. Algorytmy tworzące modele dla takich gier są często skomplikowane i obciążające obliczeniowo. Przez takie cechy dopiero od niedawna zaczęły powstawać algorytmy, zdolne pokonać ludzi w dużych grach karcianych jak "Cepheus", "DeepStack", "Libratus" lub "Pluribus". Metody zdolne tworzyć takie AI dalej są rozwijane i aktualizowane z roku na rok. W taki sposób w 2007 roku powstał CFR, potem MCCFR i

po paru latach zastąpiono je przez CFR+, aż zaczęto wykorzystywać sieci neuronowe w takich algorytmach jak Deep CFR będący tematyką pracy.

Rozdział dokładnie opisał działanie omawianego algorytmu przez przedstawienie terminów należących do działu matematyki Teoria Gier. Między innymi pokazał, że opis środowiska przez drzewa decyzyjnych pozwala na wiele uproszczeń i możliwości śledzenia gry. Dodatkowo przedstawiono problem poszukiwania stanu Równowagi Nasha, którego znalezienie jest celem większości algorytmów sztucznej inteligencji gier karcianych.

Dobrze zaimplementowany algorytm Deep CFR przy prawidłowej parametryzacji i odpowiednio dużej liczbie iteracji powinien zbliżyć się do punktu bliskiego takiego stanu. Pomimo tak optymistycznych założeń utworzenie AI potrzebuje bardzo długiego czasu nauki do znalezienia dobrych strategii gry w HULH i HUNH.

Chapter 3

Implementacja algorytmu

Program Deep CFR został napisany w niniejszej pracy, korzystając z narzędzi, pozwalających na zredukowanie nadmiarowości kodu oraz na prosta implementacje. W tym rozdziale skupiono się na opisaniu wybranych technologii, parametrów oraz funkcji, które znalazły się w implementacji algorytmu.

Baza do napisanego kodu jest język programowania, Python 3.8. Zawiera on wiele technologii wspomagających uczenie maszynowe i rozległa społeczność wspierająca jego rozwój. Poniżej wylistowano i opisano główne narzędzia użyte w pracy.

TensorFlow Jest to wysokopoziomowe API dostępne dla takich języków jak Python, JavaScript, C++ lub Java [21]. Wykorzystuje się go głównie do zadań głębokiego uczenia maszynowego. Przez swoją prostotę, dostępność i dobrą dokumentację stał się jednym z najpopularniejszych narzędzi wykorzystywanych do tworzenia sztucznych inteligencji. Dodatkowo od niedawna biblioteki technologii Keras stały się częścią Tensorflow. Daje to możliwości znacznego zredukowania kodu przy prostych problemach, które często są rozwiązane przez funkcje w tym module.

W przypadku niniejszej pracy głównie korzystano z funkcji zawartych w bibliotekach Keras. Wyjątkiem są nieliczne wiersze w kodzie gdzie np. było wymagane wykonanie obliczeń na tensorach.

Numpy Duża biblioteka do naukowych obliczeń na wielowymiarowych tablicach. Jest nieodłącznym elementem przy pisaniu programów uczenia maszynowego, zwłaszcza jeśli korzysta się z bibliotek Tensorflow. Wynika to z faktu, że wiele funkcji tego API, jako argumenty przyjmuje typy danych powiązane z Numpy [21].

Tqdm Małe narzędzie w języku Python pozwalające na wyświetlenie postępu procesów w działającym programie. Przydatne w celach testowych. W pracy zostało użyte do śledzenia iteracji drzew decyzyjnych algorytmu Deep CFR.

TensorBoard Moduł należący do API Tensorflow. Wizualizuje postępy uczenia sieci neuronowych oraz ich jakość przez przedstawienie odpowiednich wykresów.

PyPokerEngine Biblioteka wspomagająca symulacje gry Poker Texas Hold'em wraz z dodatkowymi elementami. Użyto jej jako podstawy do napisania środowiska HULH.

3.1 Implementacja sieci neuronowych

Algorytm Deep CFR do działania wymaga dwóch typów sieci neuronowych, jedna ma rozpoznawać strategie σ^t , a druga przypisana do określonego gracza przewiduje wartości D_p^t . Dodatkowo każdy z tych elementów jest trenowany na podstawie cyklicznie aktualizowanych buforów B_s i B_p . W tym rozdziale zostanie przedstawiona dokładna implementacja modeli, proces ich uczenia, struktura zbiorów danych oraz budowa środowiska, z którym jest wykonywana interakcja. W tab. 3.1 i 3.2 zamieszczono podstawowe informacje o parametrach sieci. Dalsza część rozdziału dokładnie opisuje te parametry i dodatkowe elementy, ważne podczas uczenia modelu.

Table 3.1: Podstawowe parametry sieci neuronowej θ_s .

parametr	użyte wartości
predkość uczenia	0.0001
końcowa funkcja aktywacyjna	<i>softmax</i>
rozmiar bufora	600 000
maksymalna liczba iteracji	5 000
rozmiar <i>minibatch</i>	500
parametr <i>patience</i>	40

Table 3.2: Podstawowe parametry sieci neuronowej θ_p .

parametr	użyte wartości
predkość uczenia	0.0001
końcowa funkcja aktywacyjna	<i>linear</i>
rozmiar bufora	300 000
maksymalna liczba iteracji	5 000
rozmiar <i>minibatch</i>	500
parametr <i>patience</i>	40

3.1.1 Architektura modelu

W algorytmie zaimplementowano trzy sieci neuronowe θ_1 , θ_2 , θ_s o nieskomplikowanej architekturze. Z tego powodu wykorzystano bibliotekę Keras, która do takich przypadków jest dobrym rozwiązaniem. Dodatkowo założono, że wszystkie modele będą miały podobną budowę poza ostatnią warstwą z innymi funkcjami aktywacyjnymi. Wynika to z faktu, że algorytm Deep CFR korzysta z sieci, które dostają dane wejściowe o tej samej strukturze, ale zwracają D_p^t lub σ_p^t .

Architektura sieci neuronowych składa się z 7 elementów, wejścia, wyjścia, bloku o nazwie *Flatten* [21], trzech warstw ukrytych zakończonych normalizacją i wyjściem w postaci wektora o trzech polach. Początkowo sieci dostają tablice o wymiarach (3, 52), czyli kolumny kart ze stołu, z reki oraz historie ostatniej rundy gry. Model w dalszych obliczeniach musi przekształcić takie dane do formy jedno-wymiarowej (None, 156), co robi warstwa *Flatten*. Kolejne dwa elementy składają się z 512 wag, trzecia warstwa ukryta posiada ich 256. Na trzech wymienionych elementach ustawiono funkcje *Relu*. Całość kończy się wyjściem wektora reprezentującego możliwe akcje gry HULH. Dodatkowo dane są poddawane normalizacji przez warstwę o nazwie *BatchNormalization* [21].

Wyjście modeli, aby zwracało prawidłowe liczby, używa innej funkcji aktywacyjnej niż poprzednie elementy. W przypadku predykcji strategii σ_p^t wybrano *softmax*. Sieci neuronowe θ_1 , θ_2 używają funkcji *linear*. Dokładna architektura tych obiektów jest zaprezentowana na rys. 3.1. Dodatkowo sieci korzystają z funkcji optymalizującej Adam. Predkość uczenia jest równa 0,0001, co powoduje powolne uczenie, ale minimalizuje szanse na ominięcie minimum globalnego.

Jak wynika z dokumentu prezentującego algorytm Deep CFR, uzyskuje on lepsze wyniki, jeśli sieci neuronowe są uczone za każdym razem od początku przy losowo ustawionych zerach w wagach [3]. W tym celu dodano do każdej warstwy funkcję *RandomUniform* [21], która tworzy wartości losowo w przedziałach od -0,005 do 0,005.

Ostatnim elementem sieci jest funkcja licząca błąd predykcji w trakcie uczenia. Jak wynika z opisu algorytmu Deep CFR, wymaga on zmodyfikowanej wersji MSE (*Mean Square Error*). Zostało to wykonane przy pomocy funkcji matematycznych na tensorach, jakie udostępnia Tensorflow [21]. Wszystkie te elementy zostały uwzględnione w funkcjach klasy *Poker_network* przedstawionej na rys. 3.4.

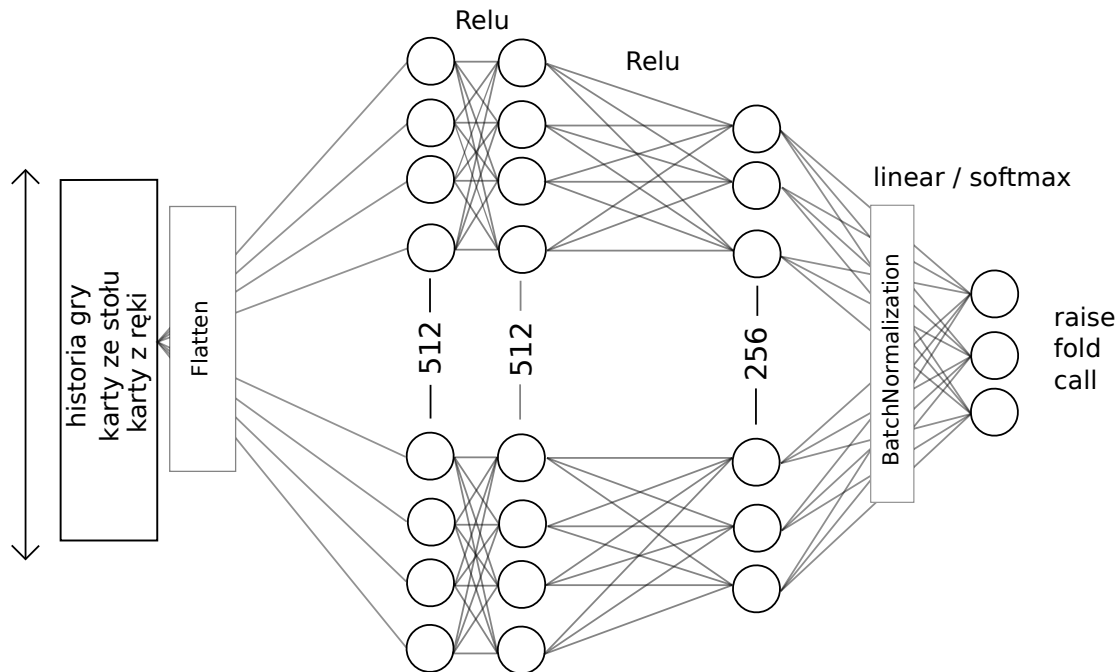


Figure 3.1: Architektura sieci neuronowych wykorzystana w programie.

3.1.2 Budowa zbiorów danych

Jak wynika z algorytmu Deep CFR, w programie muszą być zawarte trzy bufor. Maksymalna pojemność kontenerów B_1 i B_2 jest równa 300 000 próbek. Bufor B_s może posiadać ich 600 000. Każdy z dodawanych elementów do bufora składa się z trzech połączonych wektorów o rozmiarach 52. W tym celu zaimplementowano klasę *Memory* zarządzającą tymi zbiorami. Działają one jak kolejki i posiadają typ danych *deque* co oznacza, że w przypadku przepełnienia jest usuwany najstarszy wpis.

Mając uzupełnione dane w tablicach, program rozpoczyna przygotowanie zbiorów danych do uczenia wybranych sieci neuronowych. Każdy z buforów zostaje losowo przetasowany i podzielony na dwa podzbiory. Pierwszy z nich to zbiór uczący zajmujący 80% całej bazy. Wykorzystuje się go do poprawiania wag modelu sekwencyjnie. Drugim zbiorem są dane walidacyjne. Dokonanie takiego podziału było wymagane, aby przeciwdziałać stanowi przetrenowania modelu. Zbiór walidacyjny jest nadzorowany i na jego podstawie można określić moment od którego model przestaje dobrze działać. Schemat tego podziału przedstawiono na rys. 3.2. Dodatkowo w trakcie nauki sieć neuronowa aktualizuje swoje wagi w każdym kroku przez użycie elementu o nazwie *minibatch*, będącym parametrem określającym mały podzbiór bazy użyty do trenowania sieci w pojedynczym kroku. Jego liczebność jest równa 500 próbek.

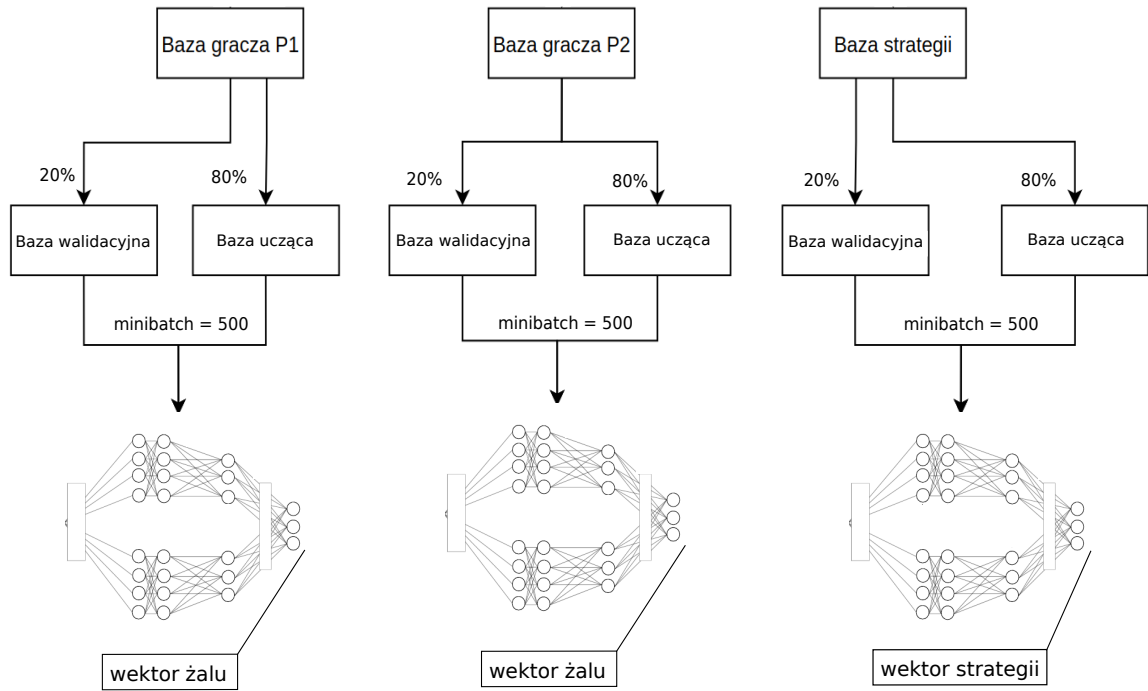


Figure 3.2: Podział danych.

3.1.3 Proces uczenia

Algorytm Deep CFR do gromadzenia obserwacji w buforach wykorzystuje metode MC-CFR ES, która eksploruje w jednej iteracji wiele razy drzewo decyzyjne. Po zakończeniu wszystkich powtórzeń zachodzi etap uczenia jednej z sieci neuronowych θ_1 , θ_2 wykonując maksymalnie 5 000 aktualizacji wag pod warunkiem jeśli wcześniej nie dojdzie do dużych oscylacji błędu predykcji w zbiorze walidacyjnym. Taki schemat działania jest wykonywany dla obu graczy i powtarza się wielokrotnie. Kończym zadaniem jest wytrenowanie sieci neuronowej θ_s . Wszystkie te elementy zostały połączone przez klasę *Brain*, agregująca obiekty *Memory* oraz *Poker_network*, rys. 3.4.

W celu śledzenia i zatrzymania oscylacji wyników uczenia, które mogły by doprowadzić model do przetrenowania, użyto funkcji *EarlyStopping* [21]. Zatrzymuje ona iteracje modelu w przypadku kiedy błąd predykcji na zbiorze walidacyjnym wzrośnie odpowiednio wysoko. Sprawdza się to na podstawie argumentu *patience*, który określa próg, po którym nauka się kończy. Taka procedura została zastosowana, aby polepszyć umiejętność generalizacji wyników przez model, a wraz z tym zwiększyć jego jakość [15]. Rys. 3.3 prezentuje przykładowy punkt zakończenia nauki modelu w momencie przekroczenia progu *patience*.

Do samego trenowania sieci neuronowych wybrano urządzenie GPU GeForce GTX 1050. Wynika to z uzyskiwania szybszych rezultatów w przeciwieństwie do innej możliwości, jaką jest wykorzystanie CPU.

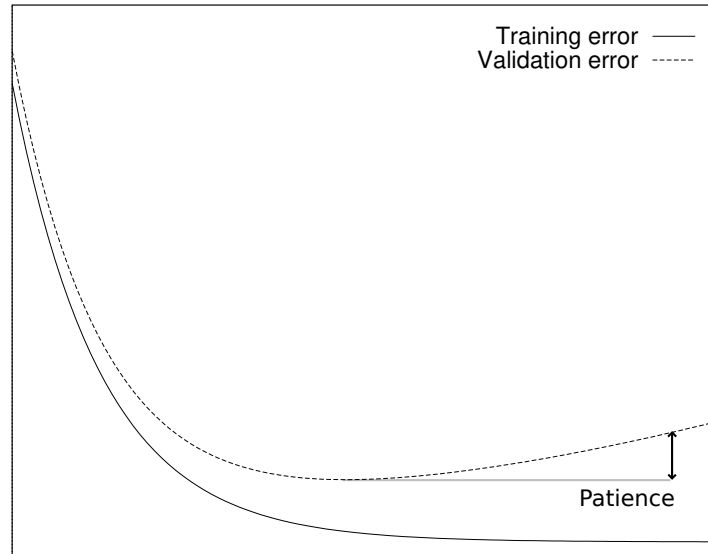


Figure 3.3: Przykład działania funkcji *EarlyStopping* [15].

3.2 Implementacja środowiska

Do symulacji gry HULH z którą model będzie się komunikował, użyto biblioteki PyPokerEngine [20]. Klasa *HULH_Emulator* zaimplementowana w programie pozwala na utworzenie obiektu zarządzającego taką grą. Przy tworzeniu go zostają zdefiniowane nazwy graczy *P1* i *P2*, które pełnią rolę identyfikatorów. W dalszej części działania algorytmu są używane do śledzenia historii gry oraz wykonywania wszelkich akcji graczy. W tabeli nr 3.3 wylistowano parametry ustawione w programie, do rozpoczęcia nauki.

Table 3.3: Parametry gry HULH.

parametr	wartość
ante	0
mała w ciemno	5
duża w ciemno	10
liczba rund po których resetuje się środowisko	1
udział każdego z graczy (stock)	80
liczba graczy	2

Zaimplementowane środowisko po każdej wprowadzonej akcji zwraca dwa elementy, obiekt *State* określający stan gry oraz historie gry w formie słownika. Obie wartości są używane do śledzenia postępu gry oraz pozycji w drzewie decyzyjnym. Dzięki takim obiektom *HULH_Emulator* może zwrócić graczowi informacje o dostępnych kartach, historii gry lub wygranych stawkach z dowolnego stanu wprowadzonego przez model.

Zaimplementowana gra symuluje HULH, czyli działa w taki sposób, że każde podbicie stawki jest wielkości dużej w ciemno. Dzięki takiemu założeniu gra nie kończy się zbyt szybko pod warunkiem, że żaden z graczy nie wykona akcji *fold*. Dodatkowo ustalono, że po każdej rundzie środowisko jest resetowane. Taki proces został ustawiony na czas uczenia. Powodem jest powstawanie mniejszych drzew decyzyjnych, co pozwala na szybszą naukę modeli oraz krótszy czas poświęcany na eksplorację gry przez MCCFR ES. W etapie rozgrywek między modelami gra będzie kończyć się dopiero po wyczerpaniu się żetonów po jednej ze stron.

3.3 Implementacja Deep CFR

Klasa *DCFR* zawiera implementację algorytmu Deep CFR. Przy tworzeniu obiektu powstają razem z nim trzy elementy klasy *Brain*, gracz, oponent oraz strategia σ . Do tego ustawiane jest środowisko *HULH_Emulator* przez referencje.

Cały proces powstawania modelu rozpoczyna się od funkcji *iterate*, która wykonuje trzy pętle *for*, dla iteracji algorytmu po 50 razy oraz powtórzeń eksploracji drzewa dla każdego z graczy po 270 razy. W pierwszej z nich środowisko tworzy nową grę, w ostatniej petli jest wykonywana funkcja *__traverse*, działająca według metody MCCFR ES. Dostaje ona argumenty *state* czyli obiekt określający stan gry, *events* - dotychczasowa historia oraz *timestep*. Ostatni argument *verbose* jest opcjonalny i służy tylko do wizualizacji powstałego drzewa decyzyjnego. Wszystko działa przez rekurencje w celu stworzenia drzewa decyzyjnego. Po zakończeniu działania MCCFR ES rozpoczyna się uczenie sieci neuronowej θ_p . Po wszystkich powtórzeniach rozpoczyna się trenowanie sieci θ_s oraz zapisanie jej w formie pliku. Aby spełnić warunek utworzenia pięciu modeli, kolejno ustawiono w funkcji *iterate* opcjonalny argument *checkpoint*. Przyjmuje on bazową wartość *None*. Przy uruchamianiu programu zmieniono ją na liczbę 10, co oznacza, że co 10 iteracji będzie trenowana sieć θ_s z uzbieranego bufora i zapisywana.

Po ustawieniu wszystkich parametrów uruchomiono program, który wykonywał się przez trzy dni. Tab. 3.4 prezentuje podstawowe parametry zaimplementowanego algorytmu.

Table 3.4: Parametry algorytmu Deep CFR.

parametr	wartość
liczba powtórzeń eksploracji drzewa	270
liczba iteracji algorytmu	50
co ile iteracji wytrenować i zapisać model	10

3.4 Podsumowanie

W tym rozdziale zaprezentowano implementację algorytmu wraz z użytymi technologiami. Program następnie został uruchomiony z zadanymi parametrami jak 50 iteracji po 270 eksploracji drzewa. Wykonywał się on przez okres czterech dni, używając GPU. Postępy tworzonych pięciu modeli były na bieżąco analizowane i zapisywane. Przedstawiono między innymi struktury sieci neuronowych, budowę zbioru danych, działanie środowiska HULH oraz funkcje związane z uczeniem algorytmu. Pokazano cel użycia takich metod jak *EarlyStopping* lub *Flatten*. Na rys. 3.4 przedstawiono dokładny schemat zaimplementowanego algorytmu. Następny rozdział przedstawi zapisane wyniki ilustrujące jakość utworzonych modeli oraz sprawdzi, który z nich będzie najlepiej grał w HULH.

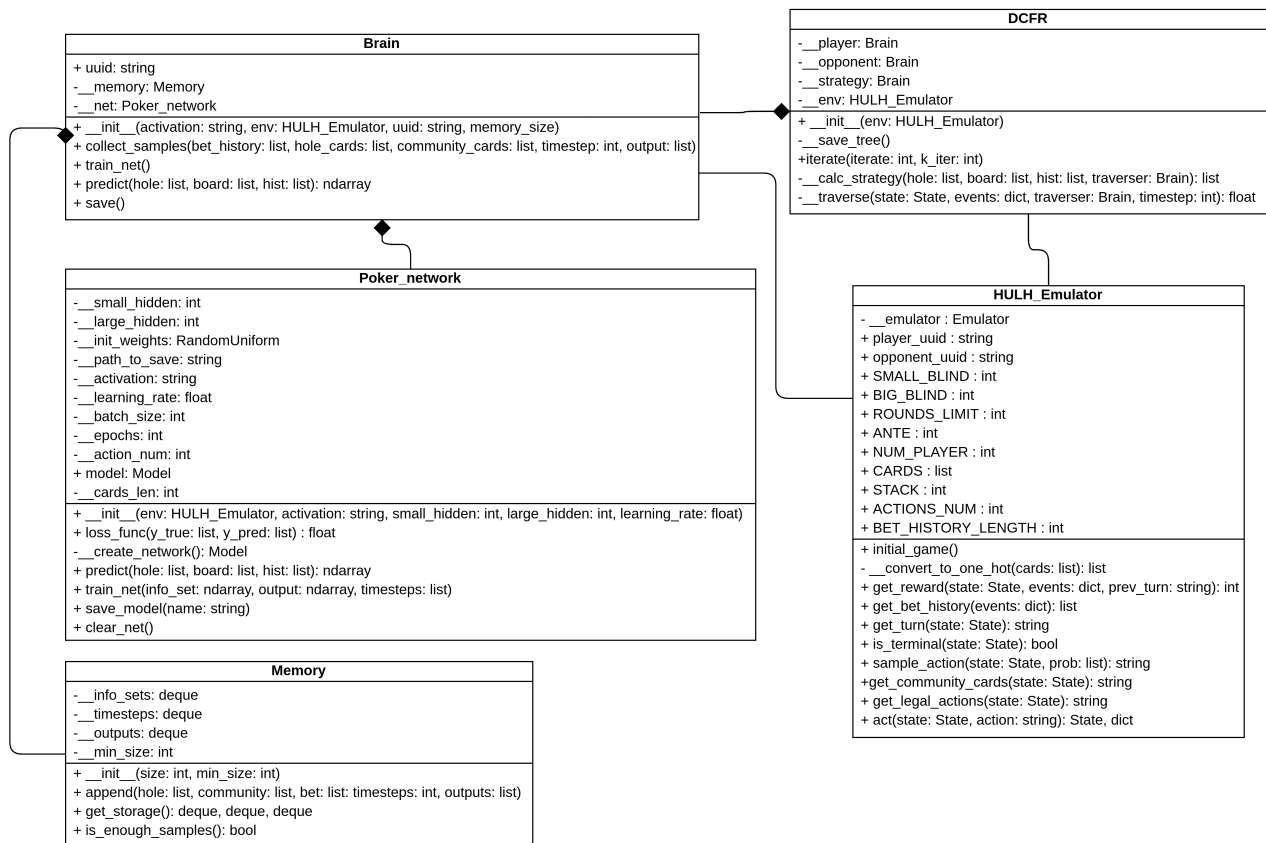


Figure 3.4: Diagram UML projektu.

Chapter 4

Wyniki

Rozdział przedstawia rezultaty powstałe po około czterech dniach nauki modelu. Celem rozdziału jest przedstawienie wyników uczenia oraz rezultatów wszystkich rozgrywek między modelami. Dodatkowo, aby przedstawić jakość algorytmu, najlepsze modele zagraja 200 razy z prostym programem zawartym w bibliotece PyPokerEngine, który symuluje gracza nigdy nieblefującego i niewykonującego akcji *raise*. Zostana tutaj zaprezentowane różnice między utworzonymi AI wraz z ich przewidywana jakością korzystając z odpowiednich wykresów bazujących na danych pobranych z modułu TensorBoard. Dodatkowo zostanie dokonana próba wyboru najlepszego modelu na podstawie wszystkich etapów oraz klasyfikacji jego stylu gry.

4.1 Proces uczenia modeli rozpoznawania

Algorytm Deep CFR zdołał utworzyć pięć modeli rozpoznawania, gdzie każdy proces nauki był śladowany przez funkcje biblioteki *TensorBoard*. Uzyskane dane następnie zostały użyte do utworzenia czytelnych wykresów ilustrujących zależność błędu predykcji od liczby iteracji.

Pierwsze utworzone AI powstał po 10 iteracjach metody Deep CFR, co zajęło około 5 godzin. Następnie rozpoczęła się nauka sieci θ_s ze zbiorem danych B_s zapełnionym przez około 300 000 próbek. Można zauważyć na rys. 4.1, że wartość błędu po 40 krokach uczenia nie zmieniła się mocno dla obu zbiorów. W przypadku danych walidacyjnych sieć neuronowa przy dokonywaniu predykcji powodowała duże oscylacje wartości, pozostając na poziomie około 2,2, co przyczyniło się do zakończenia procesu przedwcześnie przez *EarlyStopping*. Przy takim okresie model zmniejszył błąd predykcji zbioru uczącego wynoszący początkowo 0,2 do wartości bliskiej 0.

Analizując taki wykres, można stwierdzić, że model nauczył się z zebranych danych zależności między wprowadzanymi obserwacjami a zwracanymi akcjami. Dodatkowo rys. 4.6 pokazuje, że takie AI gra lepiej względem wielu utworzonych później modeli.

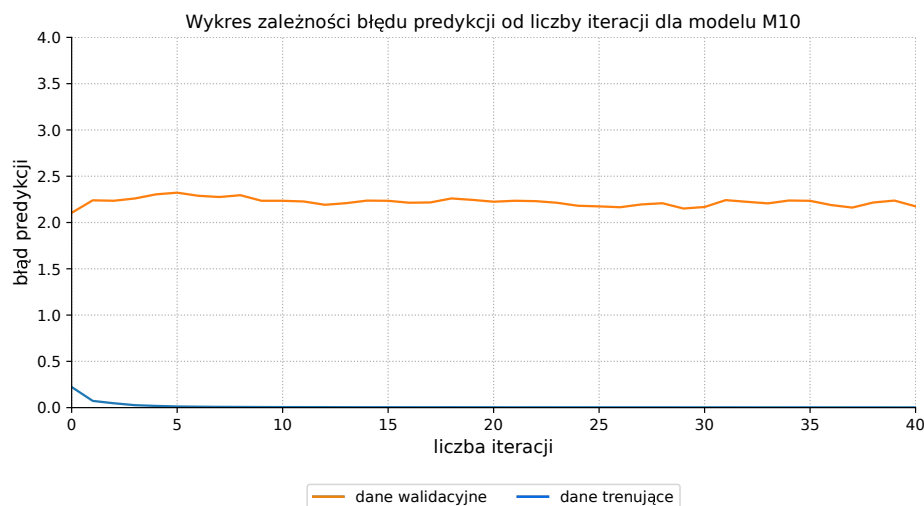


Figure 4.1: Wyniki uczenia modelu po 10 iteracjach Deep CFR.

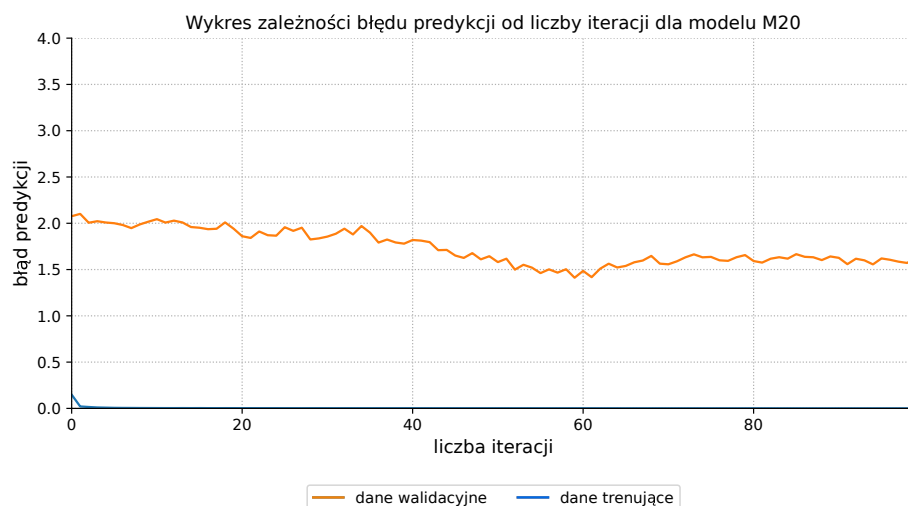


Figure 4.2: Wyniki uczenia modelu po 20 iteracjach Deep CFR.

Po 20 powtórzeniach algorytmu Deep CFR utworzono drugi model rozpoznawania. Nastąpiło to po około 12 godzinach. Algorytm do nauki użył prawie całego zapełnionego buforu B_s . W zbiorze znajdowały się nowe próbki oraz te, które zebrano przed nauką pierwszego AI. Proces uczenia przebiegał początkowo w podobny sposób jak w poprzednim etapie. Model zaczął od błędu predykcji 2,2 dla zbioru walidującego i błędu 0,2 dla buforu uczącego. Pierwsza różnica względem poprzedniego etapu jest lepsza nauka na podstawie buforu walidacyjnego, po 10 powtórzeniach wartości zaczęły mocno maleć aż do błędu predykcji wynoszącej około 1,4. Następnie pojawił się krótki wzrost wartości i utrzymywanie się na tym samym poziomie z lekkimi oscylacjami co spowodowało zakończenie procesu. Wykres 4.2 dobrze pokazuje cel używania biblioteki *EarlyStopping*. Prawdopodobnie przy większej liczbie iteracji model zacząłby uzyskiwać coraz gorsze wartości predykcji zbioru walidacyjnego. Podsumowując, model uzyskał lepsze wyniki względem etapu pierwszego, co może wynikać z większego zbioru danych. Rys. 4.6 pokazuje, że utworzony model

wygrywa nieznacznie z AI M10 przy 200 powtórzonych grach.

Kolejne AI zostało wytrenowane po 30 iteracjach. Był to drugi dzień działania algorytmu Deep CFR. W tym momencie bufor B_s był całkowicie zapełniony, przez co nowe próbki dodawane do zbioru usuwały najstarsze wpisy. Na tak dużej bazie model uczył się przez około 100 powtórzeń. Nauka modelu na podstawie zbioru uczącego była podobna jak w poprzednich etapach. Główna różnica jest bufor walidacyjny, dla którego błąd początkowo wyniósł 3,5 i zmalał do 1,7.

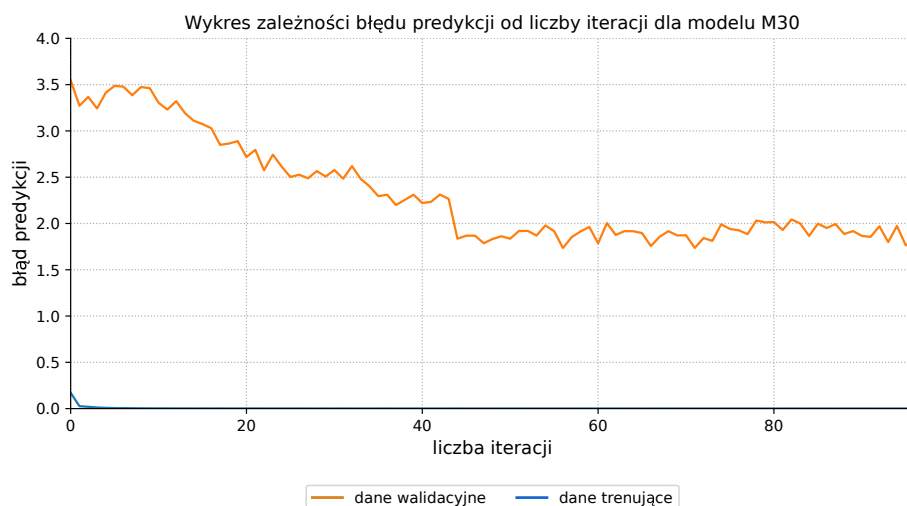


Figure 4.3: Wyniki uczenia modelu po 30 iteracjach Deep CFR.

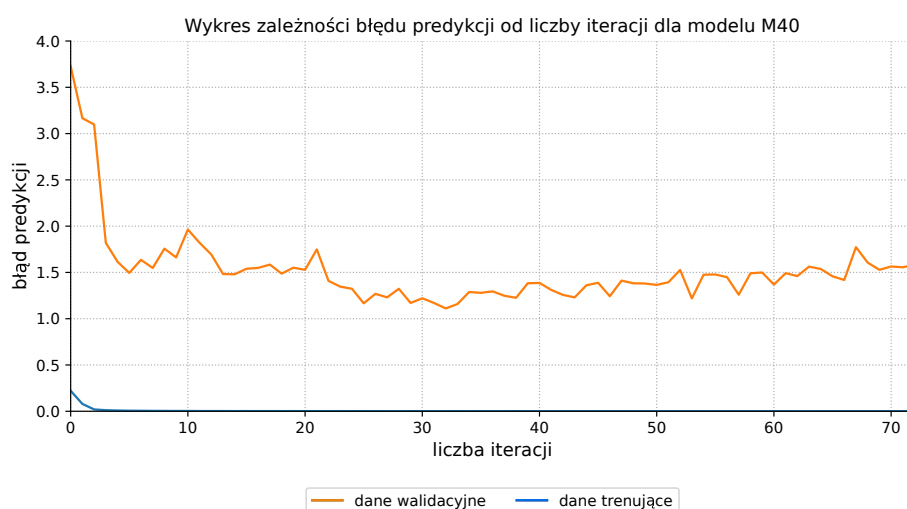


Figure 4.4: Wyniki uczenia modelu po 40 iteracjach Deep CFR.

Model nr 4 powstał po 40 iteracjach. Wytrenował sieć neuronową po około 75 powtórzeniach gdzie wartości błędu predykcji prezentują się w podobny sposób jak w poprzednim etapie. Jedyna różnica jest gwałtowny spadek błędu predykcji w pierwszych 5 iteracjach.

Ostatnie AI utworzono po 4 dniach. Rys. 4.5 pokazuje, że proces nauki przebiegał początkowo z niewielkim błędem predykcji i powoli się zmniejszał do wartości 1,9 w przypadku zbioru walidacyjnego. Drugi zbiór doprowadził do podobnych rezultatów jak we wszystkich etapach.

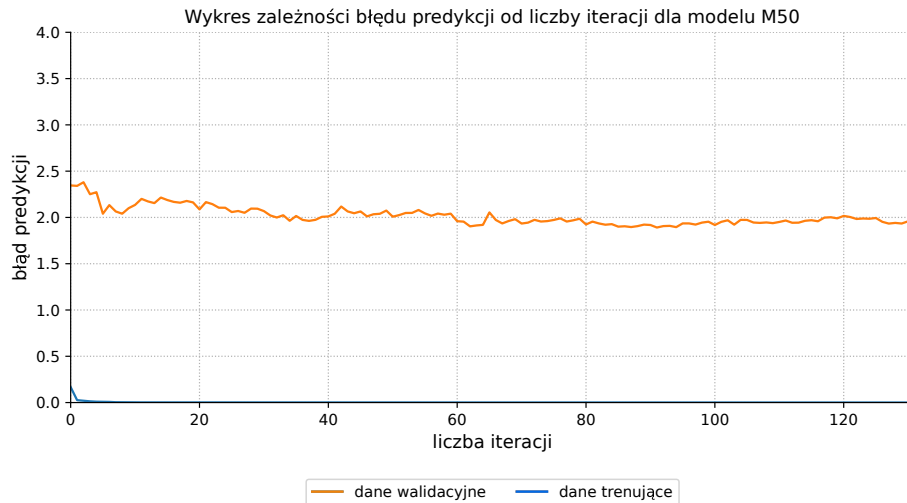


Figure 4.5: Wyniki uczenia modelu po 50 iteracjach Deep CFR.

Podsumowując powyższe wyniki można stwierdzić, że zbierane dane w buforze B_s przez Deep CFR są łatwe w znajdowaniu zależności. Główną przyczyną mogą być dane wejściowe i wyjściowe o niewielkim rozmiarze wraz z prostym środowiskiem. Prawdopodobnie przy grach dłuższych niż jedna runda, gdzie gracz przegrywa dopiero jak straci wszystkie żetony, oraz przy większym zbiorze informacji wejściowych proces uczenia by nie przebiegał tak szybko. W takim przypadku model by musiał uwzględnić dodatkowe czynniki jak numer rundy, liczba pozostałej sumy żetonów gracza lub przewaga przeciwnika.

Table 4.1: Lista utworzonych modeli.

nazwa modelu	liczba iteracji Deep CFR	liczba aktualizacji wag θ_s
M10	10	40
M20	20	98
M30	30	92
M40	40	73
M50	50	131

4.2 Wyniki rozgrywek modeli

W celu przetestowania jakości utworzonych modeli przeprowadzono 10 rozgrywek gdzie każda z nich to pojedyncza kombinacja dwóch AI. Każda gra powtórzono 200 razy z wieloma rundami tak, aby zminimalizować czynnik losowości i zakończyć grę w momencie wyczerpania się żetonów po jednej ze stron. Następnie sporządzono cztery wykresy ilustrujące wyniki tych rozgrywek. Rys. 4.6 przedstawia każdą kombinację gier z przypisaną liczbą wygranych każdemu z modeli. Rys. 4.7 i 4.8 mają za zadanie pokazać średnia wygrywanych i przegrywanych pul przez graczy, a rys. 4.9 pokazuje rozkład wykonanych akcji. Takie wykresy pozwolą stwierdzić, który z modeli częściej wygrywa, ale też częściej ryzykuje, przegrywając więcej, a który gra ostrożniej.

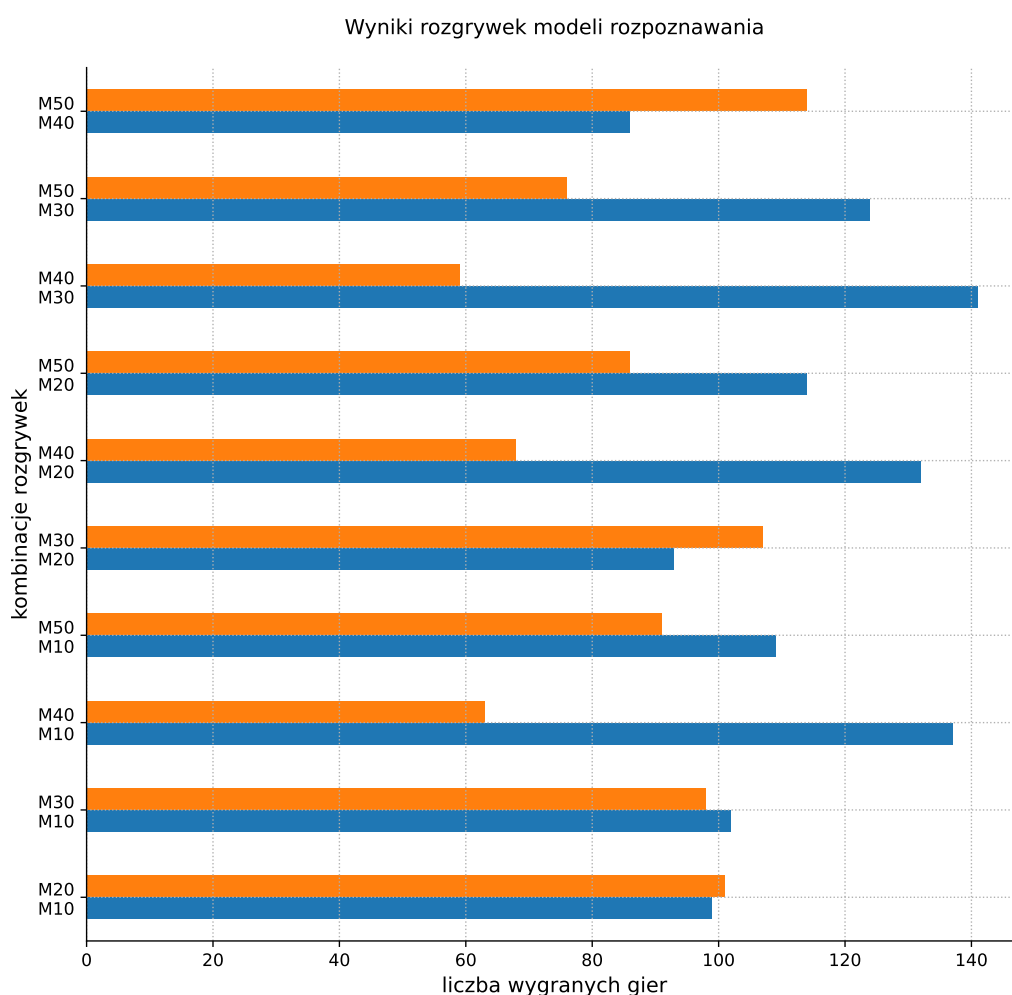


Figure 4.6: Kombinacje rozgrywek między utworzonymi modelami.

Analizując rys. 4.6 można zauważyć, że model 1 wygrywał z większością utworzonych AI najczęściej. Największą przewagę zdobył względem modelu M40, około 140 wygranych względem 60 porażek oraz około 110 zwycięstw z M50. Rozgrywki z resztą uczestników, cechuje się niewielkimi przewagami gracza. Wyjątkiem jest gracz M20, którego wybrane strategie okazały się skuteczne przeciwko oponentowi M10. Przy wielokrotnym powtarzaniu eksperymentu, modele M30 i M20 zdobywały lepsze lub gorsze wyniki niż M10 przez występujący w grze czynnik losowości. Podsumowując, gracz pierwszy nieznacznie częściej wygrywał w porównaniu M30 i M20 przez co ciężko stwierdzić czy jest od nich lepszy.

Kolejny obiekt zdobywający najlepsze wyniki to M30. Wygrał on z modelem M20 z dużą przewagą. Oznacza to, że wybierane strategie przez niego są skuteczne w rozgrywce z takim typem przeciwnika. Dodatkowo osiągnął on tak samo dużą różnicę między wygranymi a przegranymi z modelem M40. Na podstawie takich informacji można stwierdzić, że AI powstałe po 40 iteracjach przegrywa najczęściej, a utworzone w 1 i 3 etapie najrzadziej. W dalszej części rozdziału zostanie zbadana możliwa przyczyna takich wyników.

Problemem wykresu rys. 4.6 jest pokazywanie tylko częstotliwości zwycięstw modeli, ale nie zawiera informacji o sposobie gry każdego z graczy. Taka wiedza pozwoliłaby na stwierdzenie możliwych przyczyn powstałych wyników. W tym celu zebrano wszystkie wygrane wartości przez graczy z każdej rundy, a następnie obliczono ich średnie. Rys. 4.7 przedstawia uśrednione wyniki przegrywanych puli przez każdego z graczy, rys. 4.8 prezentuje odwrotną cechę. Analizując oba rysunki, można zauważyć, że model nr 1, który posiada najwięcej wygranych, zyskuje i traci największą pulę w rundach. Rys. 4.9 pokazuje też, że sztuczna inteligencja wykonuje równie często akcje *call* i *raise*, ale bardzo rzadko akcje *fold*. Jest drugim najbardziej zrównoważonym graczem, który uzyskuje przy tym najbardziej skrajne wyniki. Zaletą modelu jest to, że średnie wartości są na podobnym poziomie, co oznacza, że wygrywa i traci tyle samo. Bardzo podobnej strategii używa gracz M50, ale pomimo tego przegrywa z omawianym modelem. Gracz M30 okazał się graczem zyskującym i tracącym mniejsze wartości. Może to wynikać z częstego wykonywania akcji *call*. Modele M20, który znacznie częściej wykonuje ten ruch, zdobywa najniższą pulę wygrana i przegrana. Najgorsze wyniki ma gracz M40, który też najczęściej pasował. Uzyskał on znacznie większą średnią wartość wygranej względem przegranej, ale najrzadziej wygrywał.

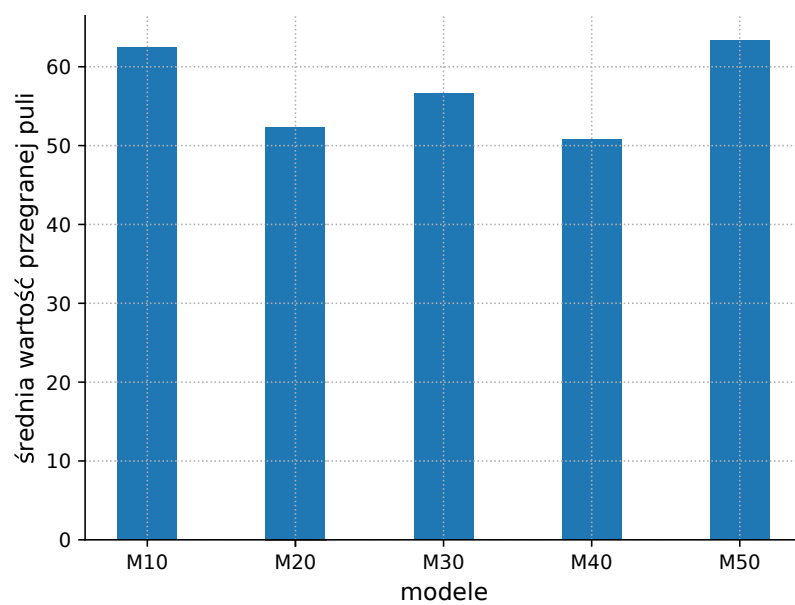


Figure 4.7: Uśrednione wartości traconych żetonów przez modele.

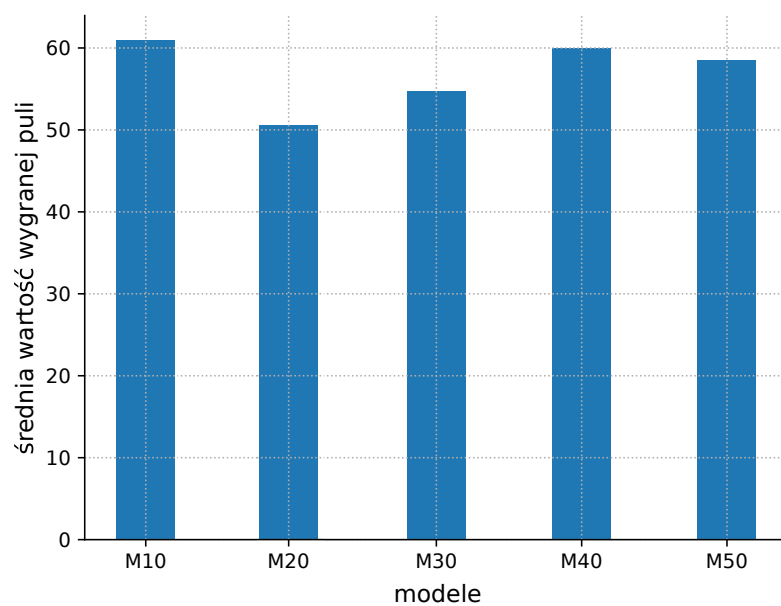


Figure 4.8: Uśrednione wartości wygrywanych żetonów przez modele.

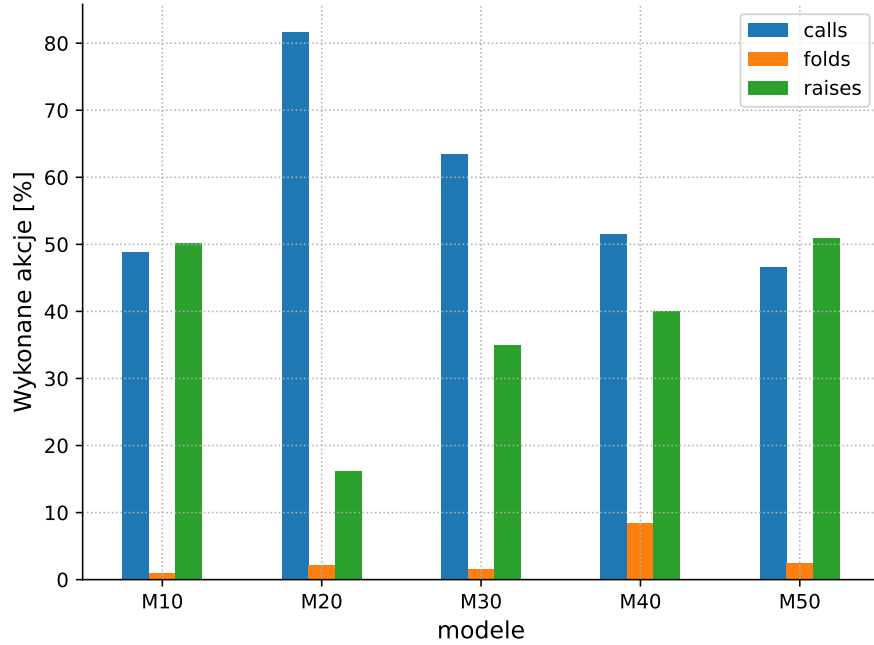


Figure 4.9: Rozkład wybieranych akcji przez modele podczas gry.

Table 4.2: Dokładne wyniki uśrednionych wartości zdobywanych lub traconych żetonów przez modele.

model	średnia wygrana	średnia przegrana
1	61	65
2	49	52
3	53	58
4	61	50
5	60	61

4.3 Porównanie modeli z programem HP

Końcowym etapem oceny jakości algorytmu Deep CFR jest wykonanie gry między najlepszymi utworzonymi modelami a graczem HP (*HonestPlayer*). Jest to program symulujący grę w HULH, wykonując tylko akcje opierające się na sile dostępnych kart. W tym celu wykorzystano funkcję *estimate_hole_card_win_rate* [20]. Wykonuje ona określona liczba iteracji możliwych wersji gier, a następnie liczy szanse wygrania z dostępnymi kartami. W zależności od zwracanej wartości wykonuje akcje *fold* lub *call*.

Rozegrano 2 gry po 200 powtórzeń, gdzie każda z nich składała się z maksymalnie 10 rund. Na rys. 4.10, 4.11, 4.12, 4.13 zaprezentowano wyniki uzyskanych rozgrywek. Pierwszy z wykresów przedstawia częstotliwość wygrywania modeli z graczem HP. Wyniki pokazują ogromną przewagę utworzonych modeli z prostym programem. Pomimo tego, analizując następne wykresy można zauważyć, że gracz HP wygrywa większą średnią

pule względem modeli i przegrywa mniejszą stawkę. Prawdopodobnie wynika to z faktu wykonywania bardzo często akcji *fold* oraz *call*. Jak widać taka strategia nie jest dobrym rozwiązaniem przez częste tracenie szansy na wygraną.

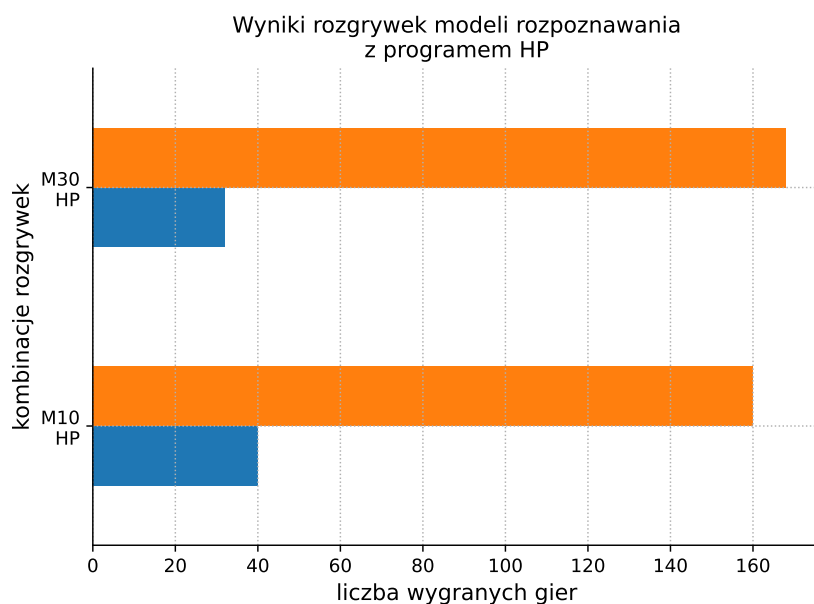


Figure 4.10: Wyniki rozgrywek między modelami M10, M30 i graczem HP.

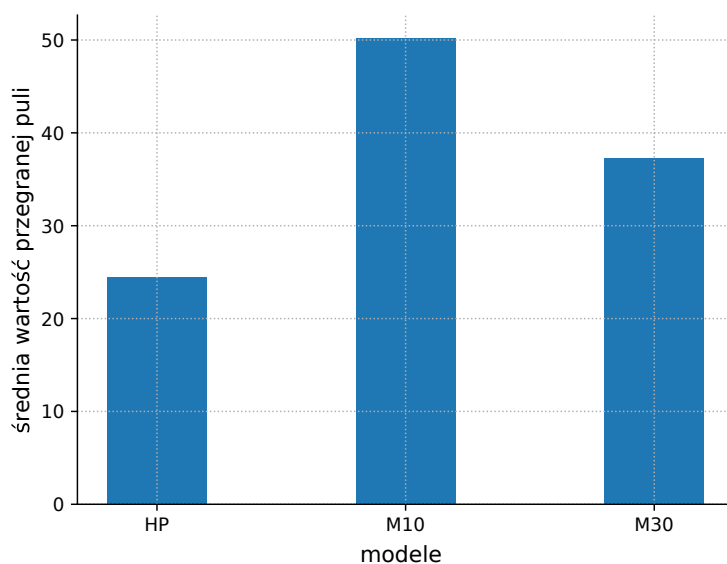


Figure 4.11: Wyniki przegrywanych pul między modelami M10, M30 i HP.

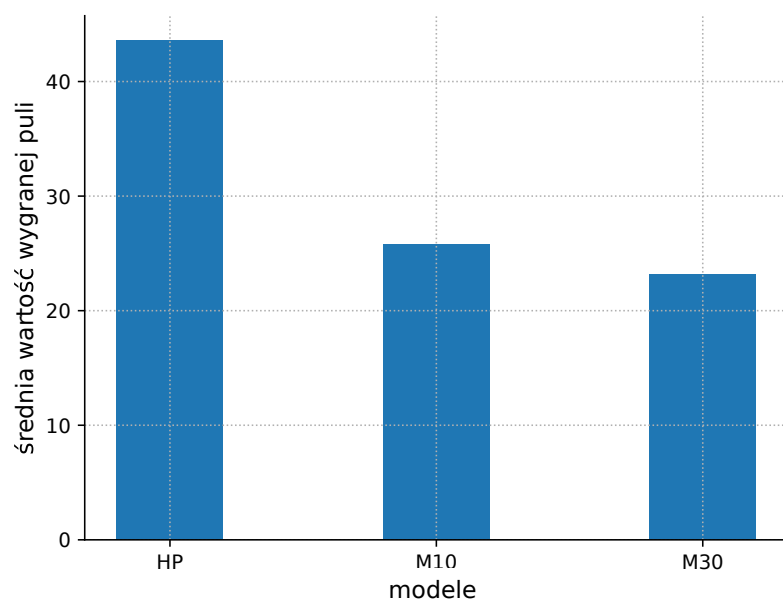


Figure 4.12: Wyniki wygrywanej pul między modelami M10, M30 i graczem HP.

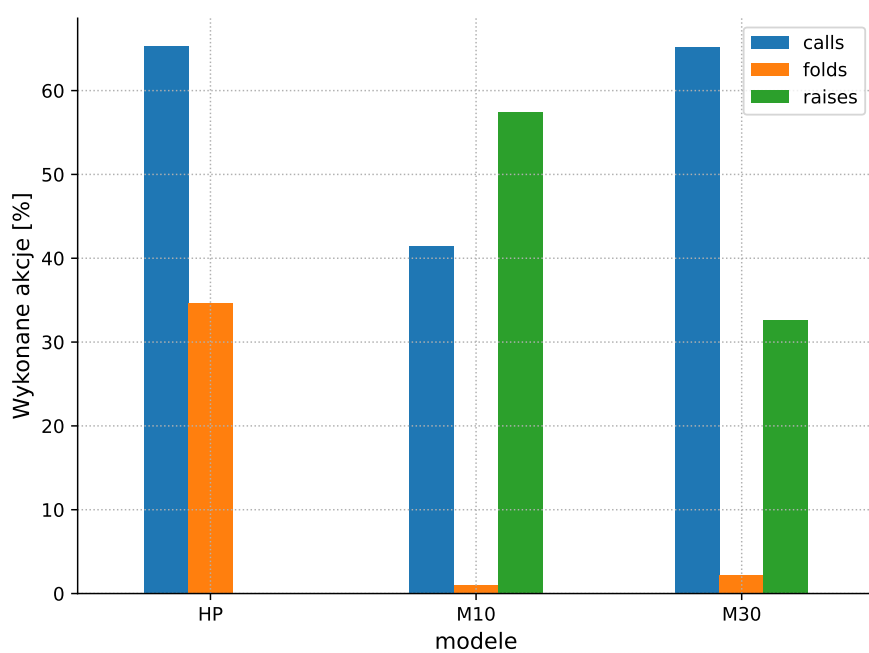


Figure 4.13: Rozkład wykonywanych akcji przez modele w trakcie rozgrywek.

4.4 Podsumowanie

Rozdział przedstawił wyniki utworzonych AI. Rezultaty okazały się mało optymistyczne. Model wytrenowany w ostatniej iteracji nie uzyskuje najlepszych wyników, a przedostatnie AI przegrywa najczęściej. Pomimo tego utworzony zbiór AI gra lepiej od prostych programów jak zaimplementowany gracz HP. Dodatkowo zauważono, że strategia polegająca na wykonywaniu tylko akcji *call* i *fold* nie daje dobrych rezultatów.

Model M10, który wygrywał najczęściej, posiada zrównoważony rozkład wykonywanych akcji *call*, *raise* oraz bardzo rzadko pasuje gre. Dodatkowo z taką strategią uzyskał najbardziej skrajne wyniki uśrednionych pul wygrywanych i przegrywanych spośród utworzonych modeli. Korzystając z informacji zawartych w rozdziale 2, można sklasyfikować graczy na bazie ich rozkładów akcji. Wszystkie AI oprócz modelu M40 wykonywały bardzo rzadko akcje *fold*, co oznacza, że wchodziły prawie zawsze do gry. Dodatkowo modele M20 i M30 wykonywały często akcje *call*, a M10 i M50 częściej *raise*. Średnia pula tracona przez te obiekty była powyżej 45 żetonów. Oznacza to, że można je prawdopodobnie sklasyfikować do grupy pomiędzy Loose Passive i Loose Aggressive.

Chapter 5

Wnioski i podsumowanie etapów pracy

Praca przedstawiła problem uczenia przez wzmacnianie w środowiskach częściowo obserwowalnych. Rozdziały 1 i 2 omówiły źródło takich problemów oraz poziom skomplikowania gry Poker Texas Hold'em. Pokazały, że należy korzystać ze skomplikowanych algorytmów, które potrafią powiązać obserwacje z najlepszymi akcjami w danym stanie. W przypadku gier karcianych w tym celu korzysta się z algorytmów bazujących na metodzie CFR.

Zaimplementowany algorytm w niniejszej pracy, Deep CFR usprawnił proces znany w metodzie CFR, przez użycie sieci neuronowych. Taka konstrukcja przyspiesza zbieranie się AI w dużych grach jak HULH do Równowagi Nasha [3]. Program użyty w pracy nie pozwolił na określenie stopnia bliskości do takiego stanu. Wynika to z faktu, że algorytmy używające metody CFR korzystają z metryki *Exploability* bazującej na wartości BR (*Best Response*) do śledzenia postępów wybieranych strategii przez AI [16]. Taki element głównie się implementuje w grach abstrakcyjnych z powodu bardzo dużej obciążalności obliczeniowej. Duże gry stosują mniej dokładną metrykę jak LBR (*Local Best Response*), która z dużym przybliżeniem zwraca stopień bliskości modelu do Równowagi Nasha [16]. Charakteryzują się one dużym skomplikowaniem implementacyjnym.

Z tego powodu do oceny jakości Deep CFR użyto prostszej metody. Rozegrano wiele gier HULH między utworzonymi modelami, a następnie przeanalizowano wyniki i wybrano najlepszego z nich na bazie uzyskanych cech. Rezultaty okazały się mało optymistyczne w porównaniu do czasu poświęconego na czas uczenia modeli. Aktualny rozdział przedstawia możliwe przyczyny takich efektów oraz wnioski po dokonaniu wszystkich poprzednich etapów pracy wraz z możliwymi poprawami parametrów. Końcowa część pracy przedstawi dalszą historię algorytmu Deep CFR. Dodatkowo zostanie przedstawiony możliwy kierunek rozwoju uczenia maszynowego w środowiskach częściowo-obserwowalnych.

5.1 Wnioski

Gra Limit Poker Texas Hold'em jest bardzo trudnym środowiskiem do uczenia sztucznych inteligencji, wynika to z częściowej obserwowalności. Z tego powodu model Cepheus, zdolny pokonywać ludzi w takim środowisku powstał dopiero w 2015 roku. Od tamtej pory rozpoczął się nagły rozwój metod uczenia maszynowego coraz większych gier karcianych. Przykładami są DeepStack, Libratus lub Pluribus rozwiązujące gre Heads Up No Limit Texas Poker Hold'em i pokonujące profesjonalnych ludzi.

Głównym zadaniem niniejszej pracy była próba zmierzenia się z takim środowiskiem. W celu uproszczenia zadania wybrano grę HULH. Następnie zaimplementowano nowoczesny algorytm Deep CFR i wytrenowano pięć modeli rozpoznawania. Pierwszy problem, jaki napotkano to, powolna nauka sieci θ_p , wraz z dużym błędem predykcji. Podczas uczenia wartości nie spadały poniżej 600. Możliwym rozwiązaniem jest zwiększenie parametru *learning rate* oraz dopracowanie architektury sieci neuronowej co przyczyni się do niewypadania wartości w minima lokalne. Dodatkowa przyczyna takich rezultatów mogą być zbiory danych o słabej jakości. Prawdopodobnie nauka by przebiegała lepiej przez zwiększenie zawartości obserwowalnych informacji wejściowych np. liczba żetonów w grze. W taki sposób sieć neuronowa by miała więcej informacji, które by zostały użyte do zwracania właściwego wektora D_p^t .

Kolejnym problemem, który może mieć duże znaczenie w grze to konstrukcja wprowadzanych informacji do modelu. Dane wejściowe użyte w implementacji to karty widziane od strony gracza oraz historia z jednej rundy. Wada takiej architektury jest to, że w praktyce rozgrywki Poker Texas Hold'em mogą odbywać się dłużej. Wtedy gracz musi dodatkowo używać takich informacji jak, wybrana strategia przeciwnika w poprzednim etapie, czy grał ostrożnie, agresywnie albo blefował. Kolejnym czynnikiem jest sposób zmieniania się gry zależnie od liczby pozostałych żetonów w puli oraz od numeru rundy. Gracze mogą podejmować bardziej ryzykowne i nierozważne ruchy, będąc w stanie bliskim porażki. Takie informacje mogłyby być kluczowe w osiągnięciu zwycięstwa. Prawdopodobnie przy uwzględnieniu tych elementów w sieci neuronowej, utworzone modele lepiej by dobierały strategię do określonych stanów gry. To by wymagało wykonywania eksploracji MCCFR ES na znacznie większych drzewach decyzyjnych obejmujących wiele rund. Dodatkowo dane wejściowe sieci neuronowej byłyby większe. W takim przypadku możliwym zbiorem informacji mógłby być zestaw składający się z widocznych kart, historii z wielu rund, liczby żetonów każdego z graczy oraz numeru gry. W taki sposób model nauczyłby się lepiej dostosowywać sposób gry do obserwacji.

Deep CFR spełnił funkcję i utworzył modele, które wygrywają z prostymi programami symulującymi grę Poker Texas Hold'em jak przedstawiony gracz nieblefujący, HP w rozdziale 4. Pomimo dobrych rezultatów dużym problemem okazał się proces powstawiania sztucznych inteligencji. Pozornie można by było oczekiwać, że algorytm będzie tworzył lepsze AI wraz z dłuższym czasem działania. W pracy doszło do odwrotnej sytuacji. Sztuczne inteligencje utworzone w iteracjach 10 i 30 wygrywały z późniejszymi obiektami. Ciężko określić przyczynę takich rezultatów. Pomocne okazałoby się użycie metryki *Exploiability* do śledzenia postępów AI pomimo zwiększenia wymaganej mocy obliczeniowej przez algorytm. Taki element pozwoliłby na stwierdzenie czy algorytm ominął punkt o najlepszej jakości i w dalszym procesie uzyskuje podobne lub coraz gorze efekty. Pozwoliło by to na zatrzymanie procesu uzyskując najlepszy możliwy model z implementacji.

5.2 Dalszy rozwój algorytmów bazujących na metodzie CFR

Algorytm powstały w 2017 roku dawał dobre rezultaty, ale przez wykorzystanie dwóch sieci neuronowych tworzył dużą wariancję wyników [17]. Z tego powodu powstał jego następca Single Deep CFR. Po wykonanych testach uzyskał nie znacznie lepsze wyniki. Algorytm dalej nie był perfekcyjny, z tego powodu w 2020 roku powstała metoda uczenia maszynowego o nazwie DREAM [18]. Jest to bardzo dobry sposób na tworzenie AI w środowiskach częściowo-obserwowalnych. Pomimo tego wada takich algorytmów jest to, że wymagają od gry, aby wartość wygranej i przegranej sumowała się do zera. Taka cecha określa się środowiska zero-sum [3]. Przez to algorytmy są mało adaptacyjne do innych środowisk. Kolejnym problemem tych metod jest przeprowadzenie testów schodzenia się do Równowagi Nasha tylko w grach 2-osobowych [3]. Wiele środowisk jak Poker Texas Hold'em standardowo uwzględnia większą ilość uczestników. Ten problem udało się rozwiązać przez model Pluribus dopiero w 2019 roku. Na podstawie takich informacji można stwierdzić, że zaimplementowany Deep CFR wraz z jego następcami nie wyczerpały tematu środowisk gier karcianych. Nawet nowsze AI, które pokonują profesjonalnych graczy jak DeepStack lub Libratus muszą trzymać się tych warunków [12] [13]. Oznacza to, że takie gry to bardzo trudne środowiska, które prawdopodobnie będą jeszcze długo badane pod względem możliwych rozwiązań.

Bibliography

- [1] Haenlein, Michael, and Andreas Kaplan. "A brief history of artificial intelligence: On the past, present, and future of artificial intelligence." *California management review* 61.4 (2019): 5-14.
- [2] Berner, Christopher, et al. "Dota 2 with large scale deep reinforcement learning." *arXiv preprint arXiv:1912.06680* (2019).
- [3] Brown, Noam, et al. "Deep counterfactual regret minimization." *International conference on machine learning*. PMLR, 2019.
- [4] Heinrich, Johannes, Marc Lanctot, and David Silver. "Fictitious self-play in extensive-form games." *International conference on machine learning*. PMLR, 2015.
- [5] Zinkevich, Martin, et al. "Regret minimization in games with incomplete information." *Advances in neural information processing systems* 20 (2007): 1729-1736.
- [6] Heinrich, Johannes, and David Silver. "Deep reinforcement learning from self-play in imperfect-information games." *arXiv preprint arXiv:1603.01121* (2016).
- [7] Teófilo, Luís Filipe Guimarães. "Building a poker playing agent based on game logs using supervised learning." (2010).
- [8] Félix, Dinis Alexandre Marialva. "Artificial intelligence techniques in games with incomplete information: opponent modelling in Texas Hold'em." (2008).
- [9] Xiang, Xuanchen, and Simon Foo. "Recent Advances in Deep Reinforcement Learning Applications for Solving Partially Observable Markov Decision Processes (POMDP) Problems: Part 1—Fundamentals and Applications in Games, Robotics and Natural Language Processing." *Machine Learning and Knowledge Extraction* 3.3 (2021): 554-581.
- [10] Nogal-Meger, P. (2012). Dylemat więźnia jako przykład wykorzystania teorii gier. *Prace i Materiały Wydziału Zarządzania Uniwersytetu Gdańskiego*, 10(4, cz. 2), 87–95.
- [11] Bowling, Michael, et al. "Heads-up limit hold'em poker is solved." *Communications of the ACM* 60.11 (2017): 81-88.
- [12] Brown, Noam, and Tuomas Sandholm. "Superhuman AI for heads-up no-limit poker: Libratus beats top professionals." *Science* 359.6374 (2018): 418-424.

- [13] Moravčík, Matej, et al. "Deepstack: Expert-level artificial intelligence in heads-up no-limit poker." *Science* 356.6337 (2017): 508-513.
- [14] Lanctot, Marc, et al. "Monte Carlo sampling for regret minimization in extensive games." *Advances in neural information processing systems* 22 (2009): 1078-1086.
- [15] Prechelt, Lutz. "Early stopping-but when?." *Neural Networks: Tricks of the trade*. Springer, Berlin, Heidelberg, 1998. 55-69.
- [16] Lisý, Viliam, and Michael Bowling. "Equilibrium approximation quality of current no-limit poker bots." *Workshops at the Thirty-First AAAI Conference on Artificial Intelligence*. 2017.
- [17] Steinberger, Eric. "Single deep counterfactual regret minimization." *arXiv preprint arXiv:1901.07621* (2019).
- [18] Steinberger, Eric, Adam Lerer, and Noam Brown. "DREAM: Deep regret minimization with advantage baselines and model-free learning." *arXiv preprint arXiv:2006.10410* (2020).
- [19] Brown, Noam, and Tuomas Sandholm. "Superhuman AI for multiplayer poker." *Science* 365.6456 (2019): 885-890.
- [20] <https://github.com/ishikota/PyPokerEngine> GitHub.GitHub repository.
- [21] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Rafal Jozefowicz, Yangqing Jia, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Mike Schuster, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. *TensorFlow: Large-scale machine learning on heterogeneous systems*, 2015. Software available from [tensorflow.org](https://www.tensorflow.org).