

Poker Project - Team 07

Arielyte Tsen Chung Ming, Devarajan Preethi, James Pang Mun Wai, Lee Yi Wei Joel, Yip Seng Yuen

National University of Singapore

chungming.tsen@u.nus.edu, e0203237@u.nus.edu, jamespang@nus.edu, lywjoel@u.nus.edu, yip@u.nus.edu

1 Introduction

The game of poker has long been a field of interest for artificial intelligence (AI) researchers. There are many challenges with developing an AI poker agent capable of competing with humans, since poker is a complex game that forces players to make decisions with incomplete and imperfect information. AI poker agents are hence faced with the task of having to make the best decision at every street under such informational constraints. Additionally, conventional poker matches impose a time limit on players to make a decision, adding another constraint for AI poker agents to work with.

This paper outlines our design of an AI poker agent for Heads Up Limit Texas Hold'em, which aims to respond competently and accurately with limited information, within a limited amount of time. In tackling the issues above, we employed a three-pronged approach:

- Using Q-learning, a reinforcement learning algorithm, to train our agent to play with the best possible strategy. (Section 3.1)
- Reduce the overall state size, by grouping similar states together. (Section 3.1, Abstraction of States)
- Precomputing the estimated strength of hand ranks, storing them in look-up tables to allow for retrieval in constant time. (Section 3.3)

The following two sections will explain our agent's strategy within each betting round, namely the pre-flop and post-flop streets.

2 Preflop Street

In order to decide what action the agent would take based only on the limited amount of information it has during the preflop street, we introduce the use of hand strength (HS), which is an estimate determined by the probability of winning with the two hole cards currently in the agent's hand. To estimate HS , we perform a Monte Carlo simulation of 10,000 games for each possible starting hand. The total number of possible starting hands is 91, which is obtained by summing up the total number of possible pocket pairs and the total number of possible hands disregarding the suits. For each game, we first pick the two starting cards for the agent, followed by randomly picking two starting cards for the opponent and the five community cards. Both players' hands are

Card 2	Card 1												
	2	3	4	5	6	7	8	9	T	J	Q	K	A
2	0.4986	0.2853	0.2921	0.3138	0.3334	0.3345	0.3702	0.3878	0.4033	0.427	0.47	0.503	0.5371
3		0.531	0.32	0.3365	0.3566	0.3624	0.3728	0.3941	0.4283	0.4583	0.4762	0.5164	0.559
4			0.5686	0.3556	0.3867	0.3834	0.3926	0.4091	0.4433	0.4712	0.504	0.5295	0.5595
5				0.5973	0.3972	0.4139	0.4093	0.4302	0.4456	0.4831	0.5121	0.5386	0.5728
6					0.6327	0.4298	0.4442	0.4552	0.4715	0.4913	0.521	0.5593	0.5883
7						0.6737	0.4608	0.4792	0.4847	0.5081	0.5341	0.5669	0.5929
8							0.6944	0.4863	0.5063	0.5276	0.5403	0.5749	0.6107
9								0.7235	0.5283	0.5423	0.5696	0.5931	0.6254
T									0.7621	0.5734	0.5924	0.6125	0.6406
J										0.7804	0.593	0.6283	0.6407
Q											0.8019	0.6368	0.6612
K												0.8373	0.6697
A													0.8556

Figure 1: Estimated winning probabilities with starting hands

then revealed and the game result is determined as a win or a loss for the agent.

We then calculate HS using the formula:

$$HS = Pr(\text{Win}) = \frac{\# \text{ of wins}}{\text{Total games played}}$$

Where the total games played is fixed at 10,000. The results of the Monte Carlo simulation are seen in Figure 1.

At every preflop street, the agent chooses an action [fold, call, raise] to take based on the HS of the starting hand:

$$\left. \begin{array}{l} \text{raise: } HS > k \\ \text{call: } HS > j \\ \text{fold: } HS \leq j \end{array} \right\} \text{ where } 0 < j < k \leq 1$$

Where j and k are pre-determined thresholds set for the agent.

3 Postflop Streets

The postflop streets consist of Flop, Turn and River. As compared to preflop street, at least three community cards are revealed, giving us more information to work with. During the postflop rounds, reinforcement learning is employed. Specifically, we use a modified version of Q-learning to train our agent.

3.1 Q-learning

In Q-learning, we maintain a table of states and actions as input for the agent. Each state has a Q-value that corresponds to each action. The Q-value signifies the expected reward at

State	Action		
	<i>fold</i>	<i>call</i>	<i>raise</i>
1			
2			
3			
...			
S			

Table 1: Q-learning Training

that state. An illustration of the Q-learning table can be seen in Table 1.

The reward function is determined based on the pot size P - the total amount that has been bet so far - at the end of the game, as seen in the formula below:

$$R_t = \begin{cases} +\frac{P}{2} & \text{if win,} \\ -\frac{P}{2} & \text{otherwise} \end{cases}$$

After each training round, the Q-value of each state is obtained by summing up the previous Q-value and the reward function for that round. More formally,

$$Q(s', a') = Q(s, a) + \sum_0^t \lambda^t R_t(s, a)$$

Where t is the number of turns between the current move and the terminal node, R_t is the reward for round t , λ is the discount factor, s is the state, a is the action and P is the pot size.

Experience Replay

Additionally, we apply the technique of experience replay to achieve more efficient use of previous experiences. We store the agent's experiences based on the formula

$$e_t = (s_t, a_t, R_t, s_{t+1})$$

The agent stores the data it discovers for each round in a table, and reinforcement learning takes place based on random sampling from the table. This reduces the amount of experience required to learn, replacing it with more computation and memory which are cheaper resources than the agent's continuous interactions with the environment [Schaul *et al.*, 2016].

ϵ -Greedy Algorithm

At every turn, with a probability of ϵ , a random action is chosen to be performed, otherwise an action with the best expected reward (highest Q-value) is chosen.

The value of ϵ will slowly decrease as the agent acquires more training. A relatively large initial value ensures that all possible paths will be explored before the agent settles into a sub-optimal pattern. The i -th *State* and possible *Actions* for the ϵ -Greedy Algorithm is

$$State_i = \{EHS_i, S_i, P_i, \#OR_i, \#SR_i, OPS_i\}$$

$$Actions = \{fold, call, raise\}$$

Where

EHS refers to the current Expected Hand Strength of a given player. More details can be found in Section 3.2.

S refers to the current Street of the round.

P refers to the Pot Size.

#OR refers to the number of Opponent Raises per street.

#SR refers to the number of Self Raises per street.

OPS refers to the Opponent Playing Style.

Abstraction of States

Without abstracting the states into groups, the total size of the state space will be very large and it would be infeasible to calculate all possible states. Therefore, we can apply abstraction by grouping some of the features in increments.

- *EHS* is grouped in increments of 0.01, from 0 to 1.
- There are three streets - Flop, Turn and River.
- *P* in a game ranges from \$0 to \$680, but by grouping it in increments of \$40, which is the 2 times big blind amount, since every raise (big blind amount - \$20) must minimally be matched by the opponent, we get a range of 0 to 17.
- *#OR* is grouped into five different groups, i.e. from 0 to 4.
- *#SR* is grouped into five different groups, i.e. from 0 to 4.
- *OPS* is grouped into four categories. More details can be found in Section 3.4.

These groupings help to reduce the complexity and cut down the size of the total state space to: $101 \times 3 \times 18 \times 5 \times 5 \times 4 = 545,400$ states

3.2 Expected Hand Strength

The Expected Hand Strength, *EHS*, is the probability of the current hand of a given player winning if the game reaches a showdown. It factors in all possible combinations of what the opponent's hand could be and the remaining hidden board cards, and compares each of these hands to the agent's hand, to see which is better. The *EHS* of the agent's hand is then calculated, by finding the number of times that particular hand turns out to be better than the opponent's. We do this by comparing the number of comparisons in which the agent had a better hand, to the total number of comparisons. More formally,

$$EHS(h) = \frac{Ahead(h) + \frac{Tied(h)}{2}}{Ahead(h) + Tied(h) + Behind(h)}$$

Where *Ahead*, *Tied*, and *Behind* functions determine respectively the number of times the player's hand wins, ties or loses the game. The above formula used to derive the *EHS* was referenced from Teofilo *et al.* [2013a].

Note that the *EHS* may be used at any round of the game. However, the time required to compute an accurate *EHS* exceeds the time constraint of 0.2s set by the project guidelines, as the number of iterations needed to compute it for a single hand at the pre-flop street is very high. This issue is addressed using look-up tables, which is explained in the next section.

3.3 EHS Look-up Table

We devised a technique similar to Average Rank Strength (ARS) [Teofilo *et al.*, 2013b] to improve the efficiency of *EHS*. In our method, three look-up tables are created - one for Flop, one for Turn and one for River. These tables consist of precomputed *EHS* values for all possible scores. The score value for a particular set of cards refers to the strength of the current hand ranking, and is computed using PyPokerEngine’s built-in *HandEvaluator.eval()* function. For example, given some combination of hole cards and community cards which contain a Pair ranking, *HandEvaluator.eval()* returns a numeric representation of this hand ranking, which would be lower than if there were a Three-of-a-Kind ranking.

Instead of computing *EHS* during gameplay, the agent can now simply refer to the look-up table, based on its current score, to obtain its *EHS*. This technique runs in constant time, and responds 1000 times faster with negligible error [Teofilo *et al.*, 2013b].

To estimate the *EHS* for the three post-flop streets, we simulated 2.5 million rounds each. Additionally, for each round of simulation, we ran the Monte Carlo sampling algorithm 500 times to obtain the *EHS*. The results are displayed below:

Street	Number of Distinct Scores
Flop	5133
Turn	13,408
River	17,470

Table 2: Result of *EHS* Look-up Tables Simulation

We are aware that not all possible scores may have been captured in the tables generated during the 2.5 million rounds. However, despite this, the results from simulating 100,000 games show that approximately 99.87 percent of the scores can be found in the table. This goes to show that our three look-up tables are highly reliable.

3.4 Opponent Playing Styles

According to Rupeneite [2010], the playing style of an opponent can be classified into four categories. Each style is distinct, in that it describes the opponent’s frequency of play and how the player bets. The four categories of playing styles are Loose/Passive, Loose/Aggressive, Tight/Passive and Tight/Aggressive. A brief description of each style is shown in Table 3 below:

Playing Styles	Description
Tight	Plays few hands and often folds.
Loose	Plays multiple and varied hands.
Aggressive	Bets and raises a lot, almost always never checking or call.
Passive	Usually checks and call, unlikely to take the lead.

Table 3: Description of Playing Styles

The Aggressive Factor, *AF*, is used to classify a player as either Aggressive or Passive. The formula for *AF* is as follows:

$$AF = \frac{\# \text{ raises}}{\# \text{ calls}}$$

Players can be classified into either Aggressive or Passive by the percentage of games they have played. Based on research by Rupeneite [2010], a threshold of 1 is used:

- Aggressive if $AF > 1$
- Passive if $AF \leq 1$

The Player Tightness, *PT*, is used to classify a player as either Loose or Tight. The formula for *PT* is as follows:

$$PT = \frac{\# \text{ folds}}{\# \text{ games}}$$

Players can be classified into either Loose or Tight by the percentage of games they have played. Based on research by Rupeneite [2010], a threshold of 0.28 is used:

- Tight if $PT < 0.28$ hands
- Loose if $PT \geq 0.28$ hands

Later, a classification process, introduced by Dinis and Reis [2008], is conducted to classify the opponent’s style of play into four categories, as seen in Table 4.

	$AF \leq 1$	$AF > 1$
$PT \geq 0.28$	Loose Passive	Loose Aggressive
$PT < 0.28$	Tight Passive	Tight Aggressive

Table 4: Style of Play Classification

4 Limitations

4.1 Q-values converging

The Q-value refers to the expected reward at every state of the game. As mentioned in the ϵ -Greedy algorithm, at every turn there is a probability of ϵ that a random action is chosen, else the action with the highest Q-value is chosen. A Q-value that is reliable is one that converges, meaning multiple occurrences of that particular state returned Q-values that did not vary too much from each other. Ideally, we would want to allow the agent to learn until all Q-values converge, which would allow the agent to choose the best action at every state. However, due to the time constraints of the project, our agent did not receive sufficient training and hence has Q-values that do not converge. This means that it is not guaranteed that our agent will always be able to take the best action, due to the inaccuracy of some of the Q-values.

4.2 Counterfactual Regret Minimisation

Counterfactual Regret Minimisation (CFR) is another reinforcement learning algorithm that we could have used while implementing our poker agent. This is an algorithm that seeks to minimise regret about its decisions at each step of a game.

There are two types of regret - positive and negative regret. Negative regret refers to the regret of having taken a particular action in a particular situation. It means the agent would have done better had it not chosen this action in this situation. Positive regret is mechanism by which the agent tracks actions that resulted in a positive outcome. Unlike Q-learning, which remembers the reward received for every action and chooses the action with the highest reward, CFR remembers the regret for every action and favours the action it regretted not having taken previously. However, CFR assumes that a terminal state is eventually reached and performs updates only after this occurs, which is not a requirement for traditional algorithms like Q-learning.

Acknowledgments

The preparation of this report would not have been possible without the help of Dr. Yair Zick and Arka Maity, National University of Singapore, School of Computing.

References

- [Dinis and Reis, 2008] Felix Dinis and Luis Paulo Reis. An experimental approach to online opponent modeling in texas hold'em poker. In *Advances in Artificial Intelligence*, pages 83–92, Savador, Brazil, October 2008. Brazilian Symposium on Artificial Intelligence.
- [Rupeneite, 2010] Annija Rupeneite. *Building Poker Agent Using Reinforcement Learning with Neural Networks*. SCITEPRESS, 2010.
- [Russell and Norvig, 2014] Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 2014.
- [Schaul *et al.*, 2016] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. In *Proceedings of the International Conference on Learning Representations 2016*, pages 1–21, San Juan, Puerto Rico, May 2016. Canada Institute for Scientific and Technical Information.
- [Teofilo *et al.*, 2013a] Luis Filipe Teofilo, Luis Paulo Reis, and Henrique Lopes Cardoso. Computing card probabilities in texas hold'em. In *Proceedings of the 8th Iberian Conference on Information Systems and Technologies*, pages 988–993, Lisbon, Portugal, June 2013. Canada Institute for Scientific and Technical Information.
- [Teofilo *et al.*, 2013b] Luis Filipe Teofilo, Luis Paulo Reis, and Henrique Lopes Cardoso. Speeding-up poker game abstraction computation: Average rank strength. In *Proceedings of the 27th Association for the Advancement of Artificial Intelligence Workshop*, pages 59–64, Bellevue, Washington, USA, July 2013. Association for the Advancement of Artificial Intelligence.