

# Project 4 - Reinforcement Learning

B4M36SMU

Monday 15<sup>th</sup> May, 2017

In your last assignment you will implement an agent capable of playing a simplified version of *blackjack* game (sometimes called *21-game*). The complete rules are in detail explained on Wikipedia [1]. However in our project we will restrict ourselves only to a simplified version.

The game is played with a standard deck of 52 cards, which is shuffled. Your goal is to score more than the dealer, however you do not want to get over 21. At the beginning you are given two cards and see one card that the dealer has. You can decide whether you draw one more card or stop playing. Once you stop playing, it is dealer's turn. The dealer has to follow a fixed strategy — as long as the sum of his cards is less than 17, he has to draw a card. Dealer stops when this condition becomes false.

Face cards (Jack, Queen and King) have value 10. Ace can be counted as 1 or 11.

At the end of the game, the player loses if the value of his cards exceeds 21. We call this situation *bust*. This holds even if the dealer busts too. If dealer busts and player not, player wins. If neither player nor dealer busts, the winner is determined by the value of cards. Who has higher sum of cards wins. Equal sums mean tie.

## 1 Implementation

We will use *OpenAI Gym* [2] library as an environment for the game. You can find the environment implementation in file `blackjack.py`. File `carddeck.py` contains a model of card, card deck and player hand. After each step your agent will get an observation as an instance of `BlackjackObservation` class and a reward. In a terminal state you get reward 1 for winning,  $-1$  for loosing and 0 for tie. In any other state you get zero reward. You are not allowed to modify files `blackjack.py` and `carddeck.py`. The same holds for file `main.py` above the comment stating that you cannot modify the code.

In file `randomagent.py` you may find a dummy agent that makes decisions completely at random. File `dealeragent.py` contains an implementation of a fixed strategy identical to strategy of the dealer. You are encouraged to check those two files and reuse the code as you want. File `tdagent.py` should contain your implementation of passive reinforcement learning agent that learns utility estimates using temporal difference. File `sarsaagent.py` should contain your implementation of SARSA. In file `evaluation.py`, you may find some ideas on how to compare various agents. You may modify the code as you want to, however it is not a requirement.

## 2 Problem Specification

1. (3 points, mandatory)  
Propose three possible nontrivial reasonable<sup>1</sup> ways how to define state in the game.
2. (1.5 points)  
For each state space representation from 1 estimate the overall number of states.

---

<sup>1</sup>For example you cannot expect points for a representation with two states - sum of values of cards is  $\leq 21$  and  $> 21$ .

You cannot just guess a number, you have to justify it somehow (e.g. by calculation). You do not need to provide an exact number, however your estimate should not be too far from the true count. If you are not sure how to calculate the number of states, you can write a program that counts them for you and submit the code together with your report.

3. (2 points, mandatory)

Pick one of the state space representations you proposed in 1. Explain why you consider it the best one and answer the following questions. Does this representation capture all information that can be used for agent decision? Or is there any simplification? If yes, will the simplification influence the result (final policy, utility values)? If yes, how much will the result be influenced? Can you use exact methods (value iteration/policy iteration) to solve the game? Why?

4. (1 point, file `tdagent.py`)

Modify your implementation of a passive reinforcement learning agent that learns utility estimates using temporal difference method. Take your implementation from the lab and make it work in the blackjack environment. Use policy that is identical to dealer's policy<sup>2</sup> and estimate value of each state.

You should have a working implementation of the agent after the lab on May 22nd. Because you will be working on the implementation in the lab, there is some cooperation allowed. Therefore the scoring for this point is low and you will get points mostly for being able to use implementation you already have. **You are not allowed to cooperate when you will be modifying your implementation to work with the blackjack environment.**

If you are not sure what you should implement, you may want to read chapter 21.2.3 in AIMA book [3] or chapter 6.1 in book [4].

5. (4 points, mandatory, file `sarsaagent.py`)

Implement SARSA algorithm.

SARSA implementation must be your own work. This means for example that if you cooperated in lab on implementation of passive reinforcement learning agent, you have to write the code again by yourself.

If you are not sure what you should implement, you may want to read chapter 21.3.2 in AIMA book [3] or chapter 6.4 in book [4].

6. (1.5 points)

Compare how successful various strategies are.

Compare random strategy (provided), dealer strategy (provided), result from 4 and the strategy learned by SARSA. You can answer for example the following questions. What is their expected or average utility? How fast do algorithms implemented in 4 and 5 learn? Does the learned utility contradict your intuition? What is utility for drawing a card when you have club nine, diamond jack and spades two in your hand and dealer has club four? What is utility of situation when you have diamond ace and spades five and dealer spades ace? Is it better to draw a card in this situation or not? Did your utility values/ $q$  values converge? Does strategy learned by SARSA follow recommendation on bottom of [1]?

---

<sup>2</sup>Draw a card if and only if sum of your cards is less than 17.

### 3 Submission and Evaluation

- All students must work individually. Cooperation on anything else than lab part of task 4 is strictly forbidden.
- Upload the results to <https://cw.felk.cvut.cz/brute/>
- Deadline is Friday June 2nd, 2017 5:00 am
- Submit all source code and a pdf report with answers to questions 1, 2 3 and 6.
- The solution must be compatible with Python version  $\geq 3.5$ .
- Project is total worth 13 points. You are required to do all parts marked as mandatory to successfully finish the project.
- Penalty for late submission is -3 points for each day of delay.
- You have to submit the assignment before going to finals. In the case that you want to take finals before the project deadline, please email me ([petr.rysavy@fel.cvut.cz](mailto:petr.rysavy@fel.cvut.cz)) and I will grade your work ASAP.
- Should you have any questions or you found a bug in code or project specification, feel free to email me and/or ask for consultation.

### References

- [1] <https://en.wikipedia.org/wiki/Blackjack>
- [2] <https://gym.openai.com/>
- [3] Russell, Stuart and Norvig, Peter. *Artificial Intelligence. "A modern approach."* Prentice-Hall, Englewood Cliffs 25 (1995): 27.  
<http://books.google.com/books?id=8jZBksh-bUMC>
- [4] Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. Vol. 1. No. 1. Cambridge: MIT press, 1998.  
<https://mitpress.mit.edu/books/reinforcement-learning>