

Modified Black Jack: A Cards Game

Instructor Mudassar Liaq

Session Fall 2024

Deadline 11:50pm on 17th July 2025

1 Objective

The goal of this assignment is to apply reinforcement learning methods to a simple card game (It's more like Blackjack). This exercise is similar to the Blackjack example in Sutton and Barto 5.3 – please note, however, that the rules of the card game are different and non-standard.

Game rules

- The game is played with an infinite deck of cards (i.e. cards are sampled with replacement).
- Each draw from the deck results in a value between 1 and 13 (uniformly distributed) with a color of red (probability 1/3) or black (probability 2/3).
- Jack is treated as 11, Queen is treated as 12, and King is treated as 13.
- At the start of the game, both the player and the dealer draw one black card (fully observed).
 - Each turn the player may either **stick** or **hit**.
 - If the player hits, then she draws another card from the deck.
 - If the player sticks, she receives no further cards.
- The values of the player's cards are added (black cards) or subtracted (red cards).
- If the player's sum exceeds 21 or becomes less than 1, then she "goes bust" and loses the game (reward -1)
- If the player sticks, then the dealer starts taking turns. The dealer always sticks on any sum of 17 or greater and hits otherwise. If the dealer goes bust, then the player wins; otherwise, the outcome – **win** (reward +1), **lose** (reward -1), or **draw** (reward 0) – is the player with the largest sum.

2 Implementation

You should write an environment that implements the game. Specifically, write a function, named *step*, which takes as input a state s (dealer's first card 1–13 and the player's sum 1–21), and an action a (hit or stick), and returns a sample of the next state s' (which may be *terminal* if the game is finished) and reward r . We will be using this environment for model-free reinforcement learning, and you should not explicitly represent the transition matrix for the MDP. There is no discounting ($\gamma = 1$). You should treat the dealer's moves as part of the environment, i.e. calling *step* with a **stick** action will play out the dealer's cards and return the final reward and terminal state.

Monte Carlo Control

Apply Monte-Carlo control to this problem. Initialize the value function to zero. Use a time-varying scalar step-size of $\alpha_t = \frac{1}{N(s_t, a_t)}$ and an ϵ -greedy exploration strategy with $\epsilon_t = \frac{N_0}{N_0 + N(s_t)}$, where $N_0 = 100$ is a constant, $N(s)$ is the number of times that state s has been visited, and $N(s, a)$ is the number of times that action a has been selected from state s . Feel free to choose an alternative value for N_0 , if it helps producing better results. Plot the optimal value function $V^*(s) = \max_a Q^*(s, a)$ using similar axes to the Fig. 1 taken from Sutton and Barto's Blackjack example.

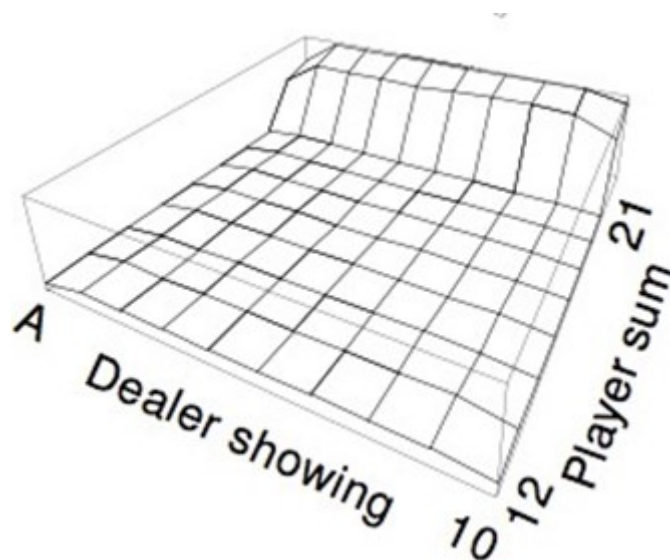


Figure 1: This is how your output will look like (except Dealer showing for you will be from 1-13 instead of 1-10)

TD Learning

Implement $Sarsa(\lambda)$. Initialize the value function to zero. Use the same step size and exploration schedules as in the previous section. Run the algorithm with parameter values $\lambda \in \{0, 0.1, 0.2, \dots, 1\}$. Stop each run after 1000 episodes and report the mean-squared error $\sum_{s,a} (Q(s,a) - Q^*(s,a))^2$ over all states s and actions a , comparing the true values $Q^*(s,a)$ computed in the previous section with the estimated values $Q(s,a)$ computed by Sarsa. Plot the mean-squared error against λ . For $\lambda = 0$ and $\lambda = 1$ only, plot the learning curve of mean-squared error against episode number.

3 Submission

- You should submit a single PDF document containing your plots and a single archive containing all your source code.
- Please organize your source code so that it is easy to follow and apparent how to run your solutions to the assignment's questions.
- You have to work on this assignment using jupyter notebooks. Any work done on Google Collab will be given 0 credit.
- You strictly have to use forward view implementation of these algorithms. Any implementations of the backward view (i.e. with eligibility traces) will result in 0 credits.