

## CS362 (Spring 2017) Machine Intelligence

### Assignment 5 – TD Learning and Monte Carlo Simulation

**Due date: 10/05/2017 (23:59)**

*This assignment has a weight of 8% on your final grade.*

The goal of this assignment is to apply reinforcement learning methods to a simple card game that we call Easy21. 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.

- 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 10 (uniformly distributed) with a colour of red (probability 1/3) or black (probability 2/3).
- There are no aces or picture (face) cards in this game
- 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.

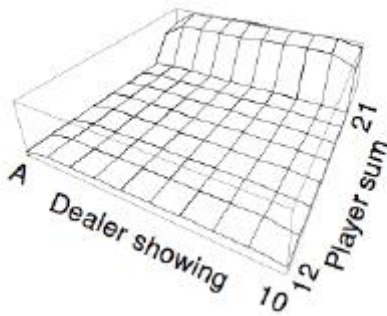
#### 1 Implementation of Easy21 (15 points)

You should write an environment that implements the game Easy21. Specifically, write a function, named **step**, which takes as input a state  $s$  (dealer's first card 1-10 and the player's sum 1-21), and an action  $a$  (hit or stick), and returns a sample of the next state  $s_0$  (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.

## 2 Monte-Carlo Control in Easy21 (25 points)

Apply Monte-Carlo control to Easy21. Initialise the value function to zero. Use a time-varying scalar step-size of  $\alpha_t = 1/N(s_t, a_t)$  and an  $\epsilon$ -greedy exploration strategy with

$\epsilon_t = N_0 / (N_0 + N(s_t))$ , 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 following figure taken from Sutton and Barto's Blackjack example.



## 3 TD Learning in Easy21 (25 points)

Implement Sarsa( $\lambda$ ) in 21s. Initialise the value function to zero. Use the same step-size and exploration schedules as in the previous section. Run the algorithm with parameter values  $\lambda = (0, 0.1, 0.2, \dots, 1)$ . 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  $\lambda$ . For  $\lambda = 0$  and  $\lambda = 1$  only, plot the learning curve of mean-squared error against episode number.

## 4 Linear Function Approximation in Easy21 (25 points)

We now consider a simple value function approximator using coarse coding. Use a binary feature vector  $\Phi(s, a)$  with  $3 \times 6 \times 2 = 36$  features. Each binary feature has a value of 1 iff  $(s, a)$  lies within the cuboid of state-space corresponding to that feature, and the action corresponding to that feature. The cuboids have the following overlapping intervals:

$$\text{dealer}(s) = \{[1, 4], [4, 7], [7, 10]\}$$

$$\text{player}(s) = \{[1, 6], [4, 9], [7, 12], [10, 15], [13, 18], [16, 21]\}$$

$$a = \{\text{hit}, \text{stick}\}$$

where

- $\text{dealer}(s)$  is the value of the dealer's first card (1-10)
- $\text{sum}(s)$  is the sum of the player's cards (1-21)

Repeat the Sarsa( $\lambda$ ) experiment from the previous section, but using linear value function approximation  $Q(s, a) = \Phi(s, a)^T \theta$ . Use a constant exploration of  $\epsilon = 0.05$  and a constant step-size of 0.01. Plot the mean-squared error against  $\lambda$ . For  $\lambda = 0$  and  $\lambda = 1$  only, plot the learning curve of mean-squared error against episode number.

## 5 Discussion (10 points)

Discuss the choice of algorithm used in the previous section.

- What are the pros and cons of bootstrapping in Easy21?
- Would you expect bootstrapping to help more in blackjack or Easy21 ? Why?
- What are the pros and cons of function approximation in Easy21?
- How would you modify the function approximator suggested in this section to get better results in Easy21?

**Prepare and upload a zip file which contains all your code and report (for the discussion section), and name this file as *<your first name>\_<your last name>\_ assignment5.zip***