

Exercise 2

Yoon-gu Hwang

August 11, 2022

1 Installation

다음과 같이 필요한 패키지를 설치합니다.

```
1 (base)$ conda create --name=rl python=3.8
2 (base)$ conda activate rl
3 (rl) $ conda install jupyter notebook matplotlib
4 (rl) $ pip install gym==0.9.6 scipy tqdm
```

2 About Environment

BlackJack 게임이 Environment로 구현되어 있습니다.

```
1 import sys
2 import gym
3 import numpy as np
4 from collections import defaultdict
5 from plot_utils import plot_blackjack_values, plot_policy
6
7 env = gym.make('Blackjack-v0')
8
9 print(env.observation_space)
10 print(env.action_space)
```

State와 Action은 다음과 같은 구조로 되어 있습니다.

```
1 Tuple(Discrete(32), Discrete(11), Discrete(2))
2 Discrete(2)
```

State와 Action에 대한 상세한 내용은 다음 docstring을 통해 자세히 알 수 있습니다.

```
1 """
2 Blackjack is a card game where the goal is to beat the dealer by obtaining cards
3 that sum to closer to 21 (without going over 21) than the dealers cards.
4 ### Description
5 Card Values:
6 - Face cards (Jack, Queen, King) have a point value of 10.
7 - Aces can either count as 11 (called a 'usable ace') or 1.
8 - Numerical cards (2-9) have a value equal to their number.
9 This game is played with an infinite deck (or with replacement).
10 The game starts with the dealer having one face up and one face down card,
11 while the player has two face up cards.
12 The player can request additional cards (hit, action=1) until they decide to stop (stick, action=0)
13 or exceed 21 (bust, immediate loss).
14 After the player sticks, the dealer reveals their facedown card, and draws
15 until their sum is 17 or greater. If the dealer goes bust, the player wins.
16 If neither the player nor the dealer busts, the outcome (win, lose, draw) is
17 decided by whose sum is closer to 21.
18 ### Action Space
```

```

19  There are two actions: stick (0), and hit (1).
20  ### Observation Space
21  The observation consists of a 3-tuple containing: the player's current sum,
22  the value of the dealer's one showing card (1-10 where 1 is ace),
23  and whether the player holds a usable ace (0 or 1).
24  This environment corresponds to the version of the blackjack problem
25  described in Example 5.1 in Reinforcement Learning: An Introduction
26  by Sutton and Barto (http://incompleteideas.net/book/the-book-2nd.html).
27  ### Rewards
28  - win game: +1
29  - lose game: -1
30  - draw game: 0
31  - win game with natural blackjack:
32    +1.5 (if natural is True)
33    +1 (if natural is False)
34  """

```

게임을 진행하는 코드는 다음과 같습니다. 매 실행마다 다른 결과가 나옵니다.

```

1  for i_episode in range(3):
2      state = env.reset()
3      while True:
4          print(state)
5          action = env.action_space.sample()
6          state, reward, done, info = env.step(action)
7          if done:
8              print('End game! Reward: ', reward)
9              print('You won :)\n') if reward > 0 else print('You lost :(\n')
10             break

```

3 Problem 1. MC Prediction

이 문제에서는 다음 `generate_episode_from_limit_stochastic()`로 Episode를 생성해주는 함수를 이용합니다.

```

1  def generate_episode_from_limit_stochastic(bj_env):
2      episode = []
3      state = bj_env.reset()
4      while True:
5          probs = [0.8, 0.2] if state[0] > 18 else [0.2, 0.8]
6          action = np.random.choice(np.arange(2), p=probs)
7          next_state, reward, done, info = bj_env.step(action)
8          episode.append((state, action, reward))
9          state = next_state
10         if done:
11             break
12     return episode

```

`generate_episode_from_limit_stochastic()`가 입력으로 들어가서 다음 함수를 채우는 것이 1번 문제입니다. First-Visit을 사용해도 되고, Every-Visit을 사용해도 됩니다.

```

1  def mc_prediction_q(env, num_episodes, generate_episode, gamma=1.0):
2      # initialize empty dictionaries of arrays
3      returns_sum = defaultdict(lambda: np.zeros(env.action_space.n))
4      N = defaultdict(lambda: np.zeros(env.action_space.n))
5      Q = defaultdict(lambda: np.zeros(env.action_space.n))
6      for i_episode in tqdm(range(1, num_episodes+1)):
7          # TODO
8          pass
9
10     return Q

```

Algorithm 1: First-Visit MC Prediction (for action values)

Input: policy π , positive integer $num_episodes$
Output: value function Q ($\approx q_\pi$ if $num_episodes$ is large enough)
Initialize $N(s, a) = 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
Initialize $returns_sum(s, a) = 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
for $i \leftarrow 1$ to $num_episodes$ do
 Generate an episode $S_0, A_0, R_1, \dots, S_T$ using π
 for $t \leftarrow 0$ to $T - 1$ do
 if (S_t, A_t) is a first visit (with return G_t) then
 $N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$
 $returns_sum(S_t, A_t) \leftarrow returns_sum(S_t, A_t) + G_t$
 end
end
 $Q(s, a) \leftarrow returns_sum(s, a) / N(s, a)$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
return Q

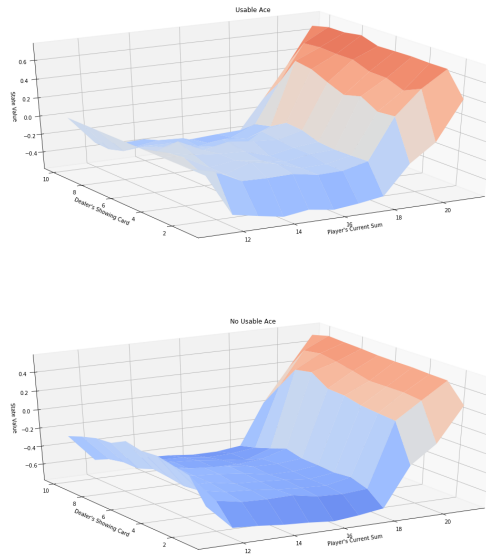


Figure 1: MC Prediction 결과

함수를 작성하고 다음 코드를 실행하면 Figure 1와 같은 그래프를 얻을 수 있습니다.

```
1 # obtain the action-value function
2 Q = mc_prediction_q(env, 500000, generate_episode_from_limit_stochastic)
3
4 # obtain the corresponding state-value function
5 V_to_plot = dict((k, (k[0] > 18) * (np.dot([0.8, 0.2], v)) + (k[0] <= 18) * (np.dot([0.2, 0.8], v)))) \
6               for k, v in Q.items())
7
8 # plot the state-value function
9 plot_blackjack_values(V_to_plot)
```

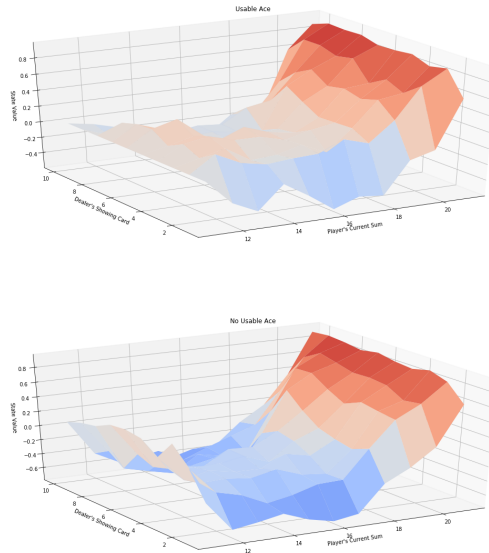


Figure 2: MC Control 결과(Value Function)

4 Problem2. MC Control

이번에는 MC Control을 구현하는 문제입니다. Episode를 생성은 ϵ -greedy 알고리즘을 통해 합니다.

Algorithm 2: First-Visit Constant- α (GLIE) MC Control

Input: positive integer $num_episodes$, small positive fraction α , GLIE $\{\epsilon_i\}$

Output: policy π ($\approx \pi_*$ if $num_episodes$ is large enough)

Initialize Q arbitrarily (e.g., $Q(s, a) = 0$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$)

```

for  $i \leftarrow 1$  to  $num\_episodes$  do
     $\epsilon \leftarrow \epsilon_i$ 
     $\pi \leftarrow \epsilon$ -greedy( $Q$ )
    Generate an episode  $S_0, A_0, R_1, \dots, S_T$  using  $\pi$ 
    for  $t \leftarrow 0$  to  $T - 1$  do
        if  $(S_t, A_t)$  is a first visit (with return  $G_t$ ) then
             $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(G_t - Q(S_t, A_t))$ 
    end
end
return  $\pi$ 

```

다음 함수 `mc_control()`의 빈칸을 채우는 것 2번째 문제입니다.

```

1 def mc_control(env, num_episodes, alpha, gamma=1.0, eps_start=1.0, eps_decay=.99999, eps_min=0.05):
2     nA = env.action_space.n
3     Q = defaultdict(lambda: np.zeros(nA))
4     epsilon = eps_start
5     # loop over episodes
6     for i_episode in tqdm(range(1, num_episodes+1)):
7         # TODO
8         pass
9
10    return policy, Q

```

코드를 다 작성한 후 다음 코드를 실행하면 Figure 2와 Figure 3과 비슷하게 나옵니다.

```

1 policy, Q = mc_control(env, 500000, 0.02)
2 V = dict((k, np.max(v)) for k, v in Q.items())
3 plot_blackjack_values(V)

```

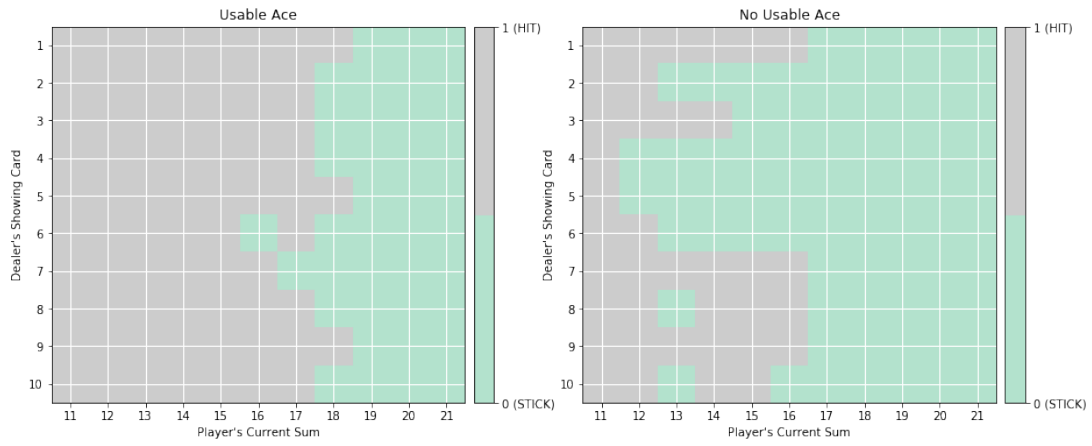


Figure 3: MC Control 결과 (Policy)

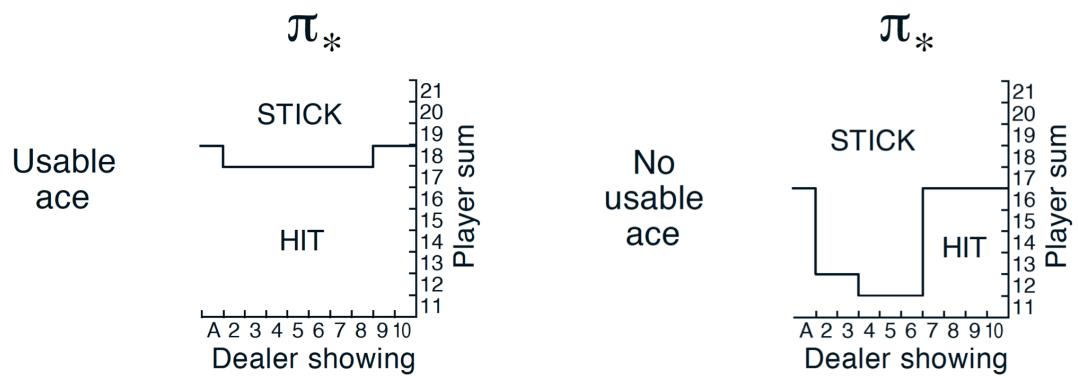


Figure 4: Optimal Policy