

Assignment 2

余北辰 519030910245

1 Experimental Purpose

- Implement first-visit MC method, every-visit MC method and TD(0) method separately
- Evaluate the uniform random policy through the three methods separately

2 Experimental Background

In this section, we present the methods of Monte-Carlo learning and Temporal-Difference learning separately.

2.1 Monte-Carlo Learning

Monte-Carlo learning is a model free method of reinforcement learning, which means learning without advance notice of MDP state transfer probability and immediate rewards. On the contrary, Monte-Carlo learning learns directly from episodes of experience. In this experiment, we use Monte-Carlo learning method to do policy evaluation.

First Visit Monte-Carlo Policy Evaluation

First visit Monte-Carlo learning use the episode under some policy π to evaluate the value function V :

$$S_1, A_1, R_1, S_2, \dots, S_T, A_T, R_T \sim \pi$$

Monte-Carlo policy evaluation uses empirical average sample returns, instead of expected return. In first visit Monte-Carlo learning, for each episode, the value function under one state is modified only when the state is first time being visited.

When a state s is firstly visited, the counter increase: $N(s) = N(s) + 1$, and the total return is increased: $S(s) = S(s) + G_t$, in which G_t is the return of the episode in time-step t :

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

Then the value function is estimated by average return $V(s) = S(s)/N(s)$. When the number of episodes is big enough, the estimated value function will convergence to the true value.

The pseudo-code for first visit Monte-Carlo policy evaluation procedure is shown below:

First-visit MC prediction, for estimating $V \approx v_\pi$

Initialize:

$\pi \leftarrow$ policy to be evaluated
 $V \leftarrow$ an arbitrary state-value function
 $Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Repeat forever:

Generate an episode using π
For each state s appearing in the episode:
 $G \leftarrow$ return following the first occurrence of s
 Append G to $Returns(s)$
 $V(s) \leftarrow \text{average}(Returns(s))$

Every Visit Monte-Carlo Policy Evaluation

Every visit Monte-Carlo policy evaluation has no difference with first visit Monte-Carlo in general, but the value function V is updated every time when a state is visited.

2.2 Temporal-Difference Policy Evaluation

Temporal-Difference learning is another model free method of reinforcement learning. Different from Monte-Carlo learning, Temporal-Difference learning learns from incomplete episodes. In this experiment we implement TD(0), which learn from just the information of the next time-step. The progress of updating the value function is shown as below:

$$V(s) = V(s) + \alpha(R + \gamma V(s') - V(s))$$

The pseudo-code for first visit Temporal-Difference policy evaluation procedure is shown below:

Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0, 1]$

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

$A \leftarrow$ action given by π for S

 Take action A , observe R, S'

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

 until S is terminal

3 Experimental Procedure

The source code for this experiment includes `gridworld.py`, `FirstVisitMC.py`, `EveryVisitMC.py`, `TD.py` and `main.py`.

3.1 gridworld.py

`gridworld.py` is exactly the same as that in the previous experiment.

3.2 FirstVisitMC.py

FirstVisitMC.py is implemented based on the above pseudo-code.

```
1 def policy_evaluation(gridworld, policy, episode_num, gamma=0.99):
2     V = np.zeros(gridworld.ns, dtype=np.float)
3     N = np.zeros(gridworld.ns, dtype=np.float)
4     for i in range(episode_num):
5         is_done = False
6         state = np.random.randint(gridworld.ns)
7         episode_reward = []
8         episode_state = []
9         while not is_done:
10             action = np.random.choice(4, 1, p=policy[state])[0]
11             prob, state, reward, is_done = gridworld.P[state][action][0]
12             episode_reward.append(reward)
13             episode_state.append(state)
14         T = len(episode_reward)
15         state_visited = []
16         for t in range(T):
17             current_state = episode_state[t]
18             if current_state not in state_visited:
19                 state_visited.append(current_state)
20                 N[current_state] += 1
21                 G = 0.0
22                 discount = 1.0
23                 for j in range(t, T):
24                     G += discount*episode_reward[j]
25                     discount *= gamma
26                 V[current_state] = V[current_state] + (G-
27                     V[current_state])/N[current_state]
28     return V
```

In each cycle, first of all, an initial state is randomly chosen. Then in each state, the action of the agent follows the policy to be evaluate. The state visited and the reward of each time-step is stored in two lists. After getting the whole episode, the value function V is updated following the equation above. It is worth noting that to guarantee "first visit", a judgment condition that whether the state is already visited is added in the algorithm.

3.3 EveryVisitMC.py

EveryVisitMC.py is modified from FirstVisitMC.py, just delete the judgment condition:

```
1 def policy_evaluation(gridworld, policy, episode_num, gamma=0.99):
2     V = np.zeros(gridworld.ns, dtype=np.float)
3     N = np.zeros(gridworld.ns, dtype=np.float)
4     for i in range(episode_num):
5         is_done = False
6         first_visit = np.zeros(gridworld.ns)
7         state = np.random.randint(gridworld.ns)
8         episode_reward = []
9         episode_state = []
10        while not is_done:
11            action = np.random.choice(4, 1, p=policy[state])[0]
12            prob, state, reward, is_done = gridworld.P[state][action][0]
13            episode_reward.append(reward)
```

```

14     episode_state.append(state)
15     T = len(episode_reward)
16     for t in range(T):
17         current_state = episode_state[t]
18         N[current_state] += 1
19         G = 0.0
20         discount = 1.0
21         for j in range(t, T):
22             G += discount*episode_reward[j]
23             discount *= gamma
24         V[current_state] = V[current_state] + (G-
V[current_state])/N[current_state]
25     return V

```

3.4 TD.py

`TD.py` is implemented based on the above pseudo-code.

```

1  def policy_evaluation(gridworld, policy, episode_num, alpha=0.5,
gamma=0.99):
2      V = np.zeros(gridworld.ns, dtype=np.float)
3      for i in range(episode_num):
4          is_done = False
5          state = np.random.randint(gridworld.ns)
6          episode_reward = []
7          episode_state = []
8          while not is_done:
9              action = np.random.choice(4, 1, p=policy[state])[0]
10             prob, state, reward, is_done = gridworld.P[state][action][0]
11             episode_reward.append(reward)
12             episode_state.append(state)
13         T = len(episode_reward)
14         for t in range(T-1):
15             current_state = episode_state[t]
16             next_state = episode_state[t+1]
17             V[current_state] = V[current_state] + alpha * \
                (episode_reward[t+1]+gamma*V[next_state]-V[current_state])
18     return V

```

The acquisition of episode is the same as MC. When updating the value function, it is just following the above equation.

3.5 main.py

`main.py` instantiates the environment, initialize the random policy and use MC or TD to evaluate that policy.

```

1  def test_first_visit_MC(gridworld, policy, episode_num):
2      V = FirstVisitMC.policy_evaluation(gridworld, policy, episode_num)
3      print('-----First Visit MC-----')
4      print('Value Function:')
5      print(np.around(np.reshape(V, gridworld.shape), decimals=4))
6      print('episode Num:')
7      print(episode_num)
8      print('-----')
9

```

```

10
11 def test_every_visit_MC(gridworld, policy, episode_num):
12     V = EveryVisitMC.policy_evaluation(gridworld, policy, episode_num)
13     print('-----Every Visit MC-----')
14     print('Value Function:')
15     print(np.around(np.reshape(V, gridworld.shape), decimals=4))
16     print('episode Num:')
17     print(episode_num)
18     print('-----')
19
20
21 def test_TD(gridworld, policy, episode_num):
22     V = TD.policy_evaluation(gridworld, policy, episode_num)
23     print('-----TD-----')
24     print('Value Function:')
25     print(np.around(np.reshape(V, gridworld.shape), decimals=4))
26     print('episode Num:')
27     print(episode_num)
28     print('-----')

```

`test_first_visit_MC()`, `test_every_visit_MC()` and `test_TD()` respectively output the results of evaluation using first visit MC, every visit MC and TD(0). The number of episodes used is also shown.

```

1 def main():
2     parser = argparse.ArgumentParser()
3     parser.add_argument('-m', '--method', type=str, default='MC1')
4     parser.add_argument('-n', '--episode_num', type=int, default=1000)
5     args = parser.parse_args()
6     method = args.method
7     episode_num = args.episode_num
8
9     gridworld = GridworldEnv([6, 6])
10    policy = np.ones([gridworld.nS, gridworld.nA]) / gridworld.nA
11
12    if method == 'MC1':
13        test_first_visit_MC(gridworld, policy, episode_num)
14
15    if method == 'MC2':
16        test_every_visit_MC(gridworld, policy, episode_num)
17
18    if method == 'TD':
19        test_TD(gridworld, policy, episode_num)

```

`main()` use parser to get the parameters entered on the command line, and depending on these parameters, the method of evaluation and the number of episodes used are chosen.

4 Experimental Result

Run first visit MC:

```
Windows PowerShell
(CS489) PS D:\Github\Reinforcement-Learning\Assignment2> python main.py -m MC1 -n 5000
-----First Visit MC-----
Value Function:
[[-14.749    0.      -20.3208 -30.3361 -35.1347 -36.5959]
 [-23.864   -21.5395 -27.7611 -32.4787 -35.6076 -36.5435]
 [-31.3801  -31.5841 -33.3157 -34.7346 -34.9922 -35.1606]
 [-36.3876  -36.656  -36.2829 -35.1606 -32.448  -30.9481]
 [-39.8033  -39.0215 -37.0623 -33.3749 -27.6801 -21.5644]
 [-40.968   -39.6457 -37.2844 -31.8397 -21.2944    0.    ]]
episode Num:
5000
(CS489) PS D:\Github\Reinforcement-Learning\Assignment2> |
```

Run every visit MC:

```
Windows PowerShell
(CS489) PS D:\Github\Reinforcement-Learning\Assignment2> python main.py -m MC2 -n 5000
-----Every Visit MC-----
Value Function:
[[-12.814    0.      -20.8159 -30.5674 -35.6855 -38.1583]
 [-23.4048  -22.3002 -28.1431 -32.8898 -36.1053 -37.71   ]
 [-31.5749  -31.6434 -33.0563 -34.776  -35.6982 -35.4033]
 [-36.3813  -36.39   -35.6837 -34.6189 -33.2997 -31.6694]
 [-39.1067  -38.3171 -36.8171 -33.572  -27.7278 -21.3383]
 [-40.1961  -38.7773 -36.5266 -31.4562 -20.556   0.    ]]
episode Num:
5000
(CS489) PS D:\Github\Reinforcement-Learning\Assignment2>
```

Run TD(0):

```
Windows PowerShell
(CS489) PS D:\Github\Reinforcement-Learning\Assignment2> python main.py -m TD -n 5000
-----TD-----
Value Function:
[[-11.3715  0.      -17.9956 -30.4545 -34.7808 -37.3711]
 [-21.0079 -21.2762 -28.3736 -32.5788 -35.2858 -36.9143]
 [-28.9879 -30.0438 -32.7403 -34.5602 -35.1826 -36.0914]
 [-36.3526 -35.3782 -36.0448 -35.2946 -33.8035 -30.9771]
 [-39.681  -38.7193 -37.2345 -33.6969 -28.1196 -23.262 ]
 [-41.3701 -40.4518 -37.5286 -31.2095 -21.5588  0.    ]]
episode Num:
5000
-----
(CS489) PS D:\Github\Reinforcement-Learning\Assignment2>
```

All the three results meet expectations, and the final evaluated value functions are similar. TD(0) runs quicker than MC methods, because that TD doesn't use the whole episode to learn, which make it quicker.

5 Experimental Summary

- Systematically learn and understand Monte-Carlo learning and Temporal-Difference learning for prediction
- Implement and successfully run the algorithm of MC and TD policy evaluation based on pseudo-code