# Reinforcement Learning Assignment 2

Chunlin Chen, 519021910463

March 22, 2023

## 1  Problem Description

| 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 6 | 7 | 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 | 16 | 17 |
| 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 |
| 30 | 31 | 32 | 33 | 34 | 35 |

Figure 1: Gridworld

Evaluate a uniform random policy: $\pi(n|\cdot) = \pi(e|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 0.25$ in the Gridworld problem illustrated in Figure 1 using the following algorithms:

- First-Visit Monte-Carlo Policy Evaluation

- Every-Visit Monte-Carlo Policy Evaluation

- TD(0) Policy Evaluation

## 2  Implementation of Model-free Prediction Algorithms

### 2.1  First-Visit Monte-Carlo Method

Let us take a look at the pseudocode of the first-visit Monte-Carlo method illustrated in Figure 2:

Here we process the episode sequence backwards temporally only for the computational convenience. The first check will make sure that the returns of

**Algorithm 1:** First-Visit MC Prediction

---

**Input**: policy $\pi$, positive integer $num\_episodes$
**Output**: value function $V$ $(\approx v_\pi,$ if $num\_episodes$ is large enough)
Initialize $N(s) = 0$ for all $s \in \mathcal{S}$
Initialize Returns$(s) = 0$ for all $s \in \mathcal{S}$
**for** $episode\ e \leftarrow 1$ **to** $e \leftarrow num\_episodes$ **do**
    Generate, using $\pi$, an episode $S_0, A_0, R_1, S_1, A_1, R_2 \ldots, S_{T-1}, A_{T-1}, R_T$
    $G \leftarrow 0$
    **for** $time\ step\ t = T - 1$ **to** $t = 0$ *(of the episode e)* **do**
        $G \leftarrow G + R_{t+1}$
        **if** $state\ S_t$ $is$ **not** $in\ the\ sequence\ S_0, S_1, \ldots, S_{t-1}$ **then**
            Returns$(S_t) \leftarrow$ Returns$(S_t) + G_t$
            $N(S_t) \leftarrow N(S_t) + 1$
    **end**
**end**
$V(s) \leftarrow \frac{\text{Returns}(s)}{N(s)}$ for all $s \in \mathcal{S}$
**return** $V$

---

Figure 2: The Pseudocode of First-Visit Monte-Carlo

a certain state will only be updated at its first appearence in one episode. The code implementation is shown below:

```python
def firstVisitMonteCarlo(env, policy, num_episodes, discount_factor=1.0):
    N = {j:0 for j in range(env.nS)}
    V = np.zeros(env.nS)

    for episode_index in range(1, num_episodes + 1):
        if episode_index % 1000 == 0:
            print("\rEpisode {}/{}.".format(episode_index, num_episodes),
                end="")
            sys.stdout.flush()
        # Generate an episode, which is an array of (state, action,
            reward) tuples
        episode = []
        state = env.reset()
        G = 0

        for step in range(100):
            action = random.choice(policy[state])
            next_state, reward, done = env.step(state, action)

            actions = ['N', 'E', 'S', 'W']
            action = actions[action]

            episode.append((state, action, reward))
            if done:
```

```python
            break
        state = next_state

    for index, state_tuple in enumerate(episode[::-1]):
        G = discount_factor * G + state_tuple[2]
        # first visit check
        if str(state_tuple[0]) not in np.asarray(episode)[:,
             0][:len(episode)-index-1]:
            N[state_tuple[0]] += 1
            # incremental update
            V[state_tuple[0]] = V[state_tuple[0]] + (G -
                 V[state_tuple[0]]) / N[state_tuple[0]]

return V
```

You may notice that we actually applied the incremental Monte-Carlo updates, therefore, the variable Returns were no longer used, we updated the V(s) directly instead.

## 2.2   Every-Visit Monte-Carlo Method

The pseudocode of the every-visit Monte-Carlo method is shown in Figure 3. We can see that the only difference of every-visit Monte-Carlo from first-visit

---

**Algorithm 2:** Every-Visit MC Prediction

**Input**: policy $\pi$, positive integer $num\_episodes$
**Output**: value function $V$ ($\approx v_\pi$, if $num\_episodes$ is large enough)
Initialize $N(s) = 0$ for all $s \in \mathcal{S}$
Initialize Returns$(s) = 0$ for all $s \in \mathcal{S}$
**for** $episode$ $e \leftarrow 1$ **to** $e \leftarrow num\_episodes$ **do**
  Generate, using $\pi$, an episode $S_0, A_0, R_1, S_1, A_1, R_2 \ldots, S_{T-1}, A_{T-1}, R_T$
  $G \leftarrow 0$
  **for** $time\ step\ t = T - 1$ **to** $t = 0$ $(of\ the\ episode\ e)$ **do**
    $G \leftarrow G + R_{t+1}$
    Returns$(S_t) \leftarrow$ Returns$(S_t) + G_t$
    $N(S_t) \leftarrow N(S_t) + 1$
  **end**
**end**
$V(s) \leftarrow \frac{\text{Returns}(s)}{N(s)}$ for all $s \in \mathcal{S}$
**return** $V$

---

Figure 3: The Pseudocode of Every-Visit Monte-Carlo

Monte-Carlo is that the former does not do the first check and it updates the returns of a certain state as long as it appears in the episode sequence. Therefore, every-visit Monte-Carlo can be easily implemented by slightly modifying the implementation of first-visit Monte-Carlo. And the code is here:

```python
def everyVisitMonteCarlo(env, policy, num_episodes, discount_factor):
    N = {j:0 for j in range(env.nS)}
    V = np.zeros(env.nS)

    for episode_index in range(1, num_episodes + 1):
        if episode_index % 1000 == 0:
            print("\rEpisode {}/{}.".format(episode_index, num_episodes),
                end="")
            sys.stdout.flush()
        # Generate an episode, which is an array of (state, action,
            reward) tuples
        episode = []
        state = env.reset()
        G = 0

        for step in range(100):
            action = random.choice(policy[state])
            next_state, reward, done = env.step(state, action)

            actions = ['N', 'E', 'S', 'W']
            action = actions[action]

            episode.append((state, action, reward))
            if done:
                break
            state = next_state

        for index, state_tuple in enumerate(episode[::-1]):
            G = discount_factor * G + state_tuple[2]

            N[state_tuple[0]] += 1
            # incremental update
            V[state_tuple[0]] = V[state_tuple[0]] + (G -
                V[state_tuple[0]]) / N[state_tuple[0]]

    return V
```

## 2.3   TD(0) Policy Evaluation

The pseudocode of the TD(0) method is shown in Figure 4.
The implementation is here:

```python
def TemporalDifference0(env, policy, num_episodes, discount_factor,
    learning_rate):
    V = np.zeros(env.nS)

    for episode_index in range(1, num_episodes + 1):
        if episode_index % 1000 == 0:
```

```
        print("\rEpisode {}/{}.".format(episode_index, num_episodes),
              end="")
        sys.stdout.flush()

    state = env.reset()

    while True:
        action = random.choice(policy[state])
        next_state, reward, done = env.step(state, action)

        V[state] = V[state] + learning_rate * (reward +
            discount_factor * V[next_state] - V[state])
        state = next_state

        if done:
            break

return V
```

---

**Tabular TD(0) for estimating $v_\pi$**

Input: the policy $\pi$ to be evaluated
Algorithm parameter: step size $\alpha \in (0, 1]$
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop for each episode:
  Initialize $S$
  Loop for each step of episode:
    $A \leftarrow$ action given by $\pi$ for $S$
    Take action $A$, observe $R, S'$
    $V(S) \leftarrow V(S) + \alpha \big[ R + \gamma V(S') - V(S) \big]$
    $S \leftarrow S'$
  until $S$ is terminal

Figure 4: The Pseudocode of TD(0)

# 3 Results Analysis

## 3.1 Monte-Carlo Prediction

### 3.1.1 First-Visit MC

We set the discount factor $\gamma = 0.9$ and the number of episodes be 10000. The final value function under the uniform random policy is shown in Figure 5:

Figure 5: The Result of First-Visit Monte-Carlo

### 3.1.2 Every-Visit MC

We set the discount factor $\gamma = 0.9$ and the number of episodes be 10000. The final value function under the uniform random policy is shown in Figure 6:



Figure 6: The Result of Every-Visit Monte-Carlo

## 3.2 TD(0) Prediction

We set the discount factor $\gamma = 0.9$, the learning rate $\alpha = 1^{-5}$, and the number of episodes be 500000. The final value function under the uniform random policy is shown in Figure 7:



Figure 7: The Result of TD(0)

## 3.3 Conclusion

In the experiment, we find that compared to the MC method, TD(0) requires much more episodes to converge, but it also runs much faster.

6