

# 强化学习 实验一

---

1.实验目的

2.实验原理

3.实验内容

Exp1 Monte-Carlo

Exp2 Sarsa

Exp3 Q-learning

姓名：吕奇澹

学号：SA23006171

## 1.实验目的

1. 理解、学习蒙特卡罗算法原理，并编码实现first-visit版本
2. 理解、学习TD算法原理，并编码实现SARSA、Q-learning

## 2.实验原理

特卡洛算法是一类通过模拟随机样本来解决问题的方法，通常可分为两类：

1. **问题具有内在的随机性：** 对于问题本身包含随机性的情况，利用计算机运算能力直接模拟这种随机过程。
2. **问题可转化为随机分布的特征数：** 将问题转化为某种随机分布的特征数，例如随机事件的概率或随机变量的期望值。通过随机抽样方法，估计随机事件发生的频率或估算随机变量的数字特征，并将其作为问题的解。

**Sarsa算法：** Sarsa算法是一种on-policy方法，其特点是原始策略和更新策略一致。不同于Monte Carlo方法的地方在于，Sarsa算法在执行完一个动作后就可以更新其值函数，而不需要采样一个完整的轨迹。这种方法适用于需要在每个时间步骤都更新值函数的情况。

**Q-learning算法：**Q-learning算法是一种off-policy方法。其原始策略和值函数更新策略不一致，而且不需要采样整个轨迹进行策略更新。不同于Sarsa算法的是，Q-learning在更新值函数时使用的是贪心策略而不是e-greedy策略。这意味着在更新值函数时，Q-learning选择当前状态下最有价值的动作，而不是基于某个概率选择其他动作。

## 3.实验内容

### Exp1 Monte-Carlo

Step1: Generate an episode:

```
Python |
1  # step 1 : Generate an episode.
2      # An episode is an array of (state, action, reward) tuples
3  def generate_one_episode(env, generate_policy):
4      trajectory = []
5      state = env.reset()
6      for i in range(1000):
7          Pi_table = generate_policy(state)
8          action = np.random.choice(np.arange(len(Pi_table)), p=Pi_table)
9          next_state, reward, done, _ = env.step(action)
10         trajectory.append((next_state, action, reward))
11         if done:
12             break
13         state = next_state
14     return trajectory
15     trajectory = generate_one_episode(env, policy)
```

Step 2 : Find all (state, action) pairs we've visited in this episode:

```
Python |
1  s_a_pairs = set([(x[0], x[1]) for x in trajectory])
```

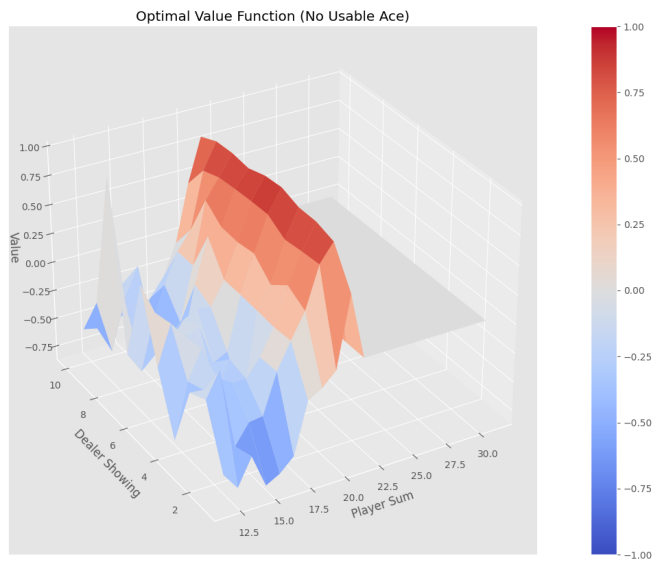
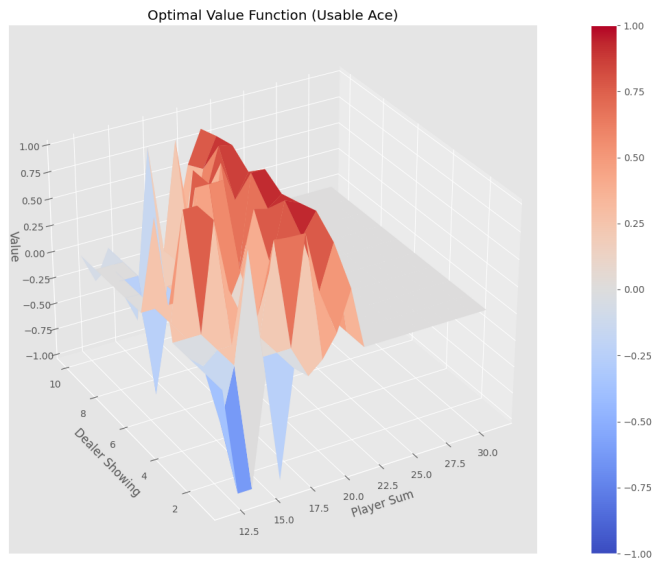
Step 3 : Calculate average return for this state over all sampled episodes:

- First-visit version:

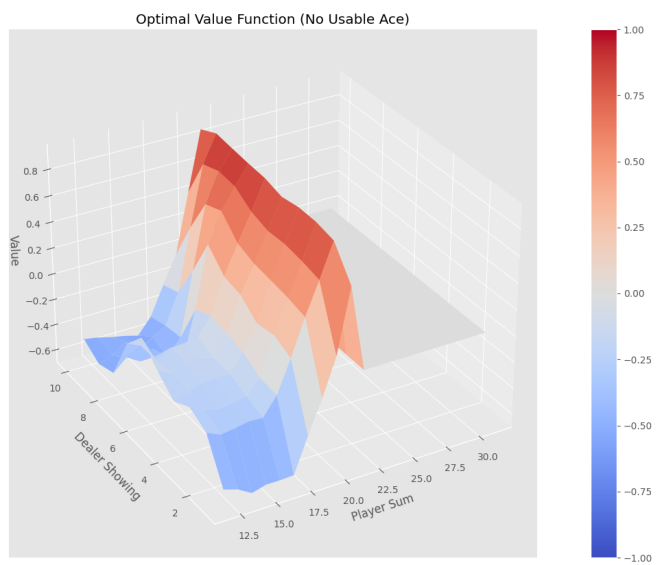
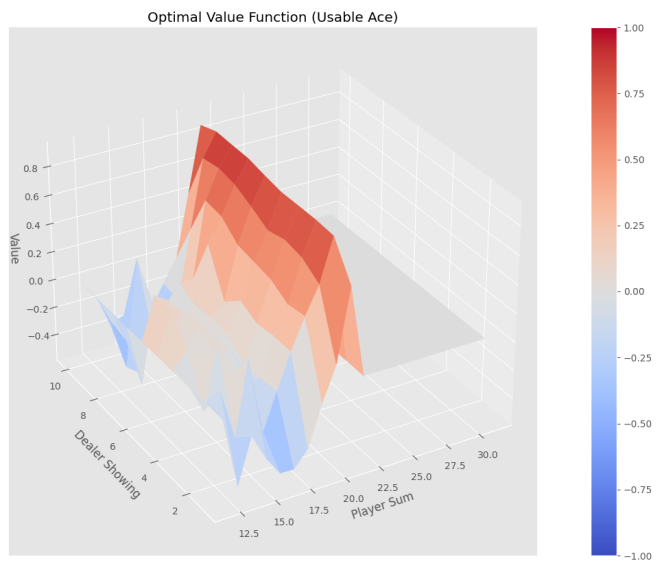
```
1 for state, action in s_a_pairs:
2     s_a = (state, action)
3     first_visit_id = next(i for i, x in enumerate(trajecory) if x
4     [0] == state and x[1] == action)
5     G = sum([x[2] * (discount_factor ** i) for i, x in enumerate(tr
6     ajectory[first_visit_id:])])
7     returns_sum[s_a] += G
8     returns_count[s_a] += 1.
9     Q[state][action] = returns_sum[s_a] / returns_count[s_a]
```

- Corresponding Results:

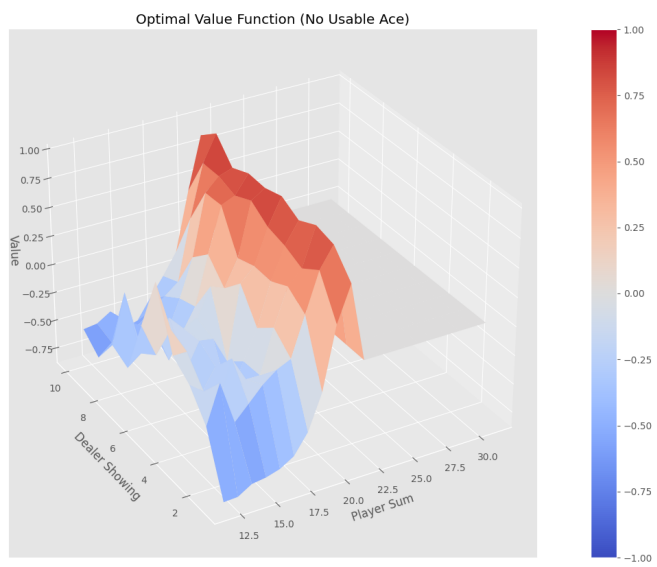
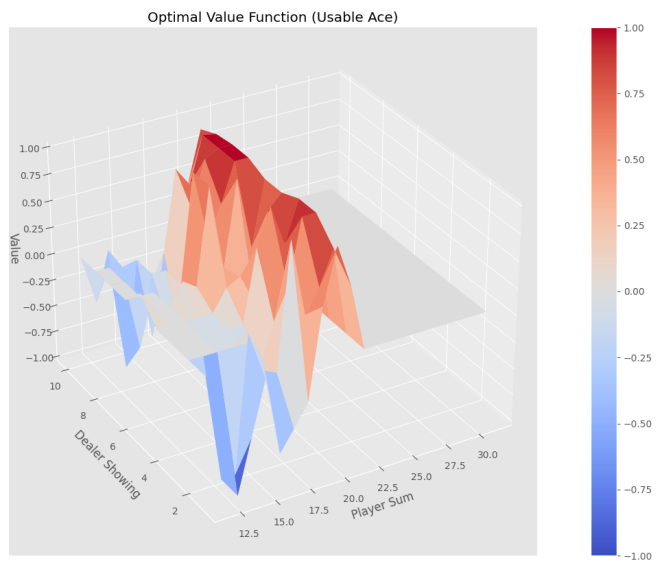
1. 10000 episodes:



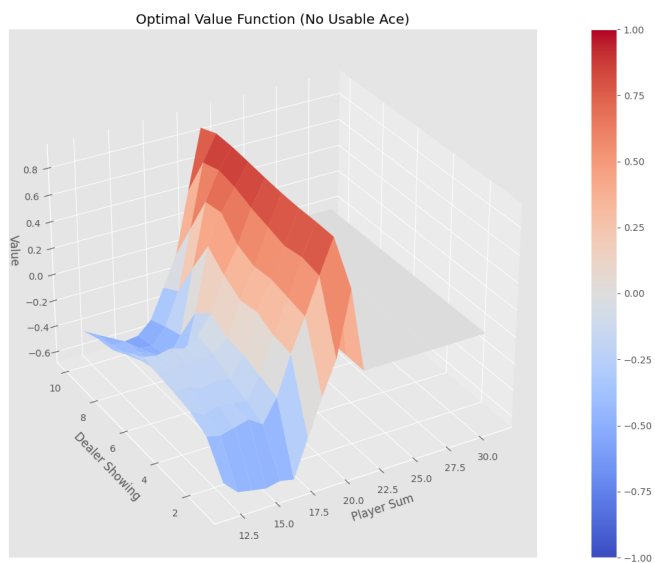
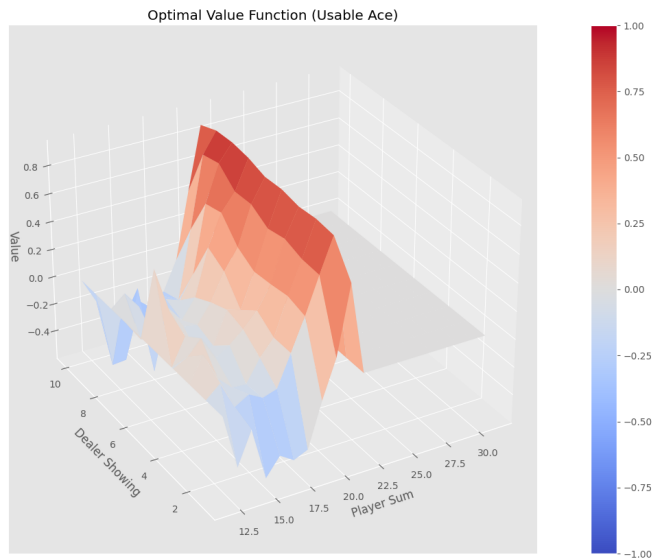
2. 500000 episodes:



- Every-visit version:
1. 10000 episodes:



2. 500000 episodes:



- 对比every-visit和first-visit，可以发现它们的主要区别在于对重复出现的相同状态 计算对应的reward.对于first-visit，每种状态仅计算第一次出现的reward；对于every-visit，每种状态计算reward均值。

## Exp2 Sarsa

step 1 : Take a step(1 line code, tips : env.step()):

```
1 new_state, reward, done, _ = env.step(action)
```

Step 2: Choose the next action:

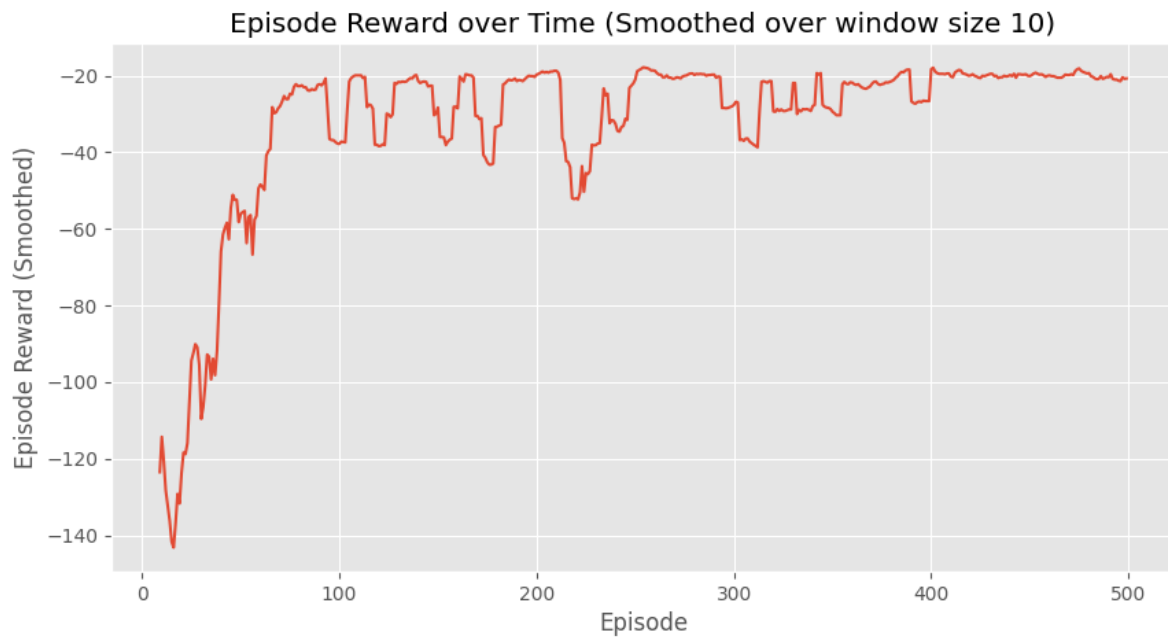
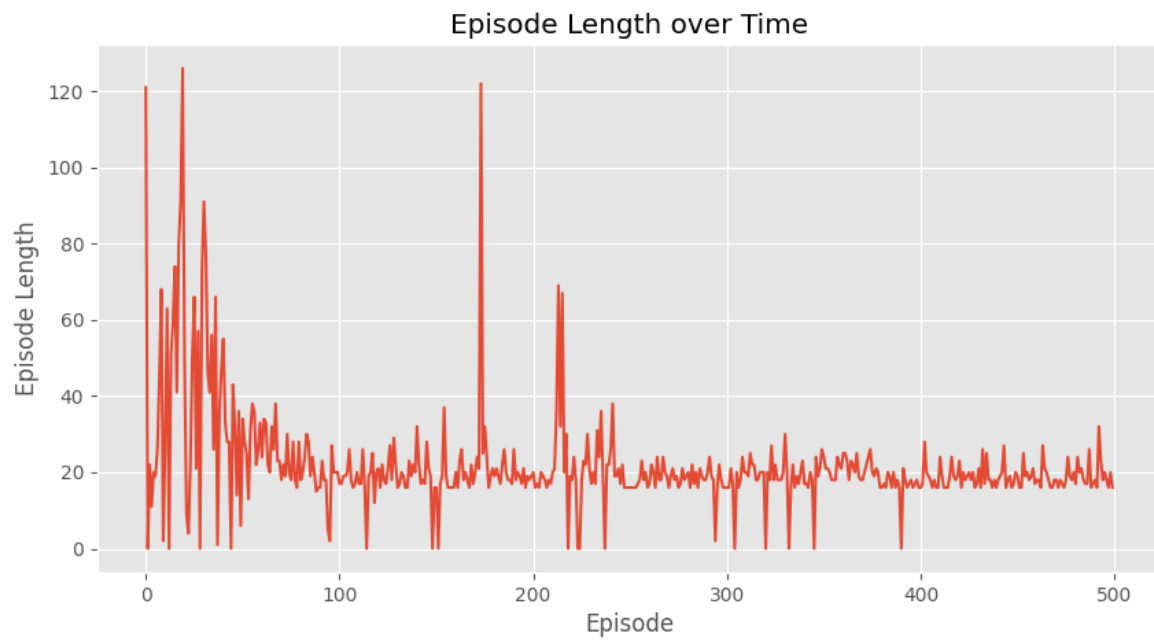
```
1 new_action = np.random.choice(np.arange(env.action_space.n), p=policy(new_s
tate))
2 stats.episode_rewards[i_episode] += reward
3 stats.episode_lengths[i_episode] = t
```

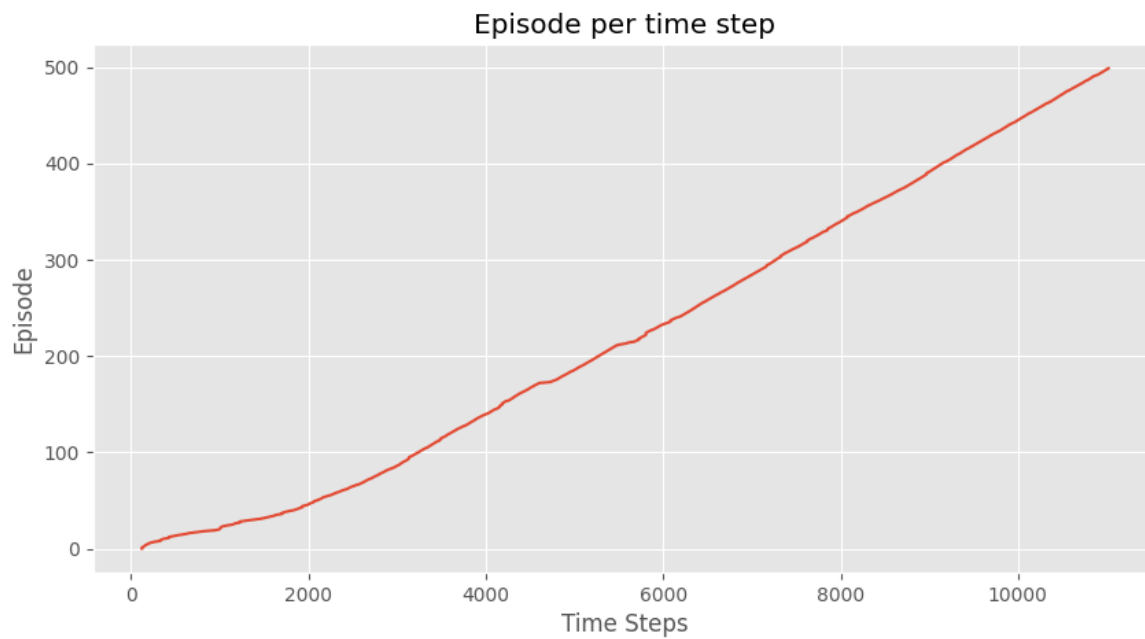
step 3 : Pick the next action

```
1 if done:
2     Q[state][action] = Q[state][action] + alpha * (reward + discount_factor
r * 0.0 - Q[state][action])
3     break
4 else:
5     Q[state][action] = Q[state][action] + alpha * (reward + discount_factor
r * Q[new_state][new_action] - Q[state][action])
6     state = new_state
7     action = new_action
```

- Corresponding Results:







## Exp3 Q-learning

### 1. 实现Q-learning 算法

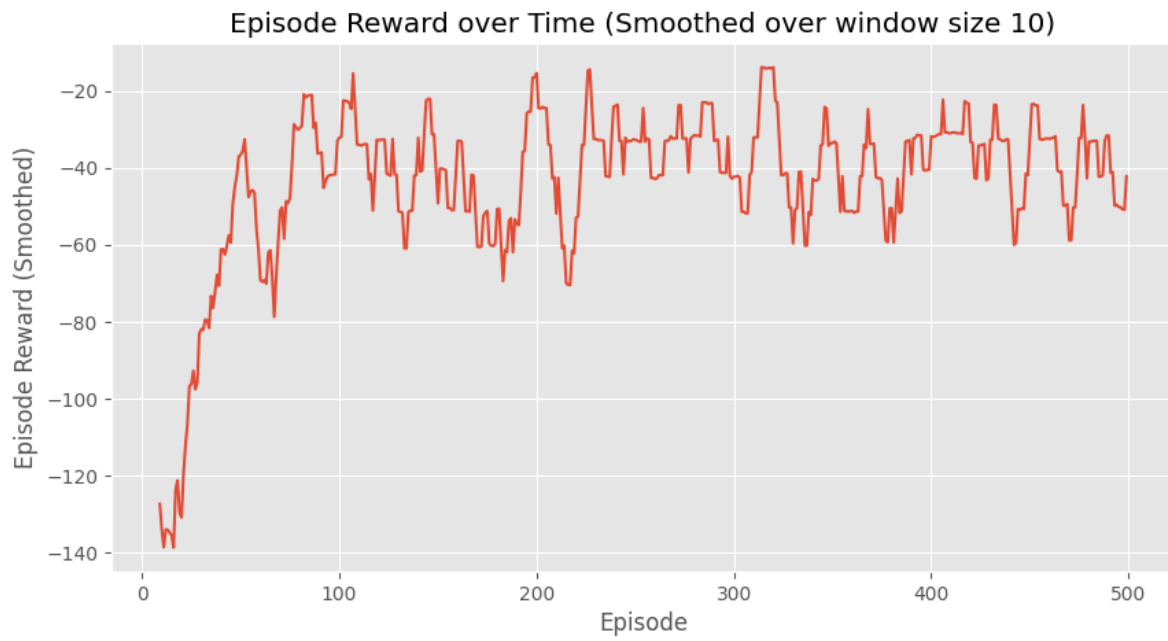
step 1 : Take a step:

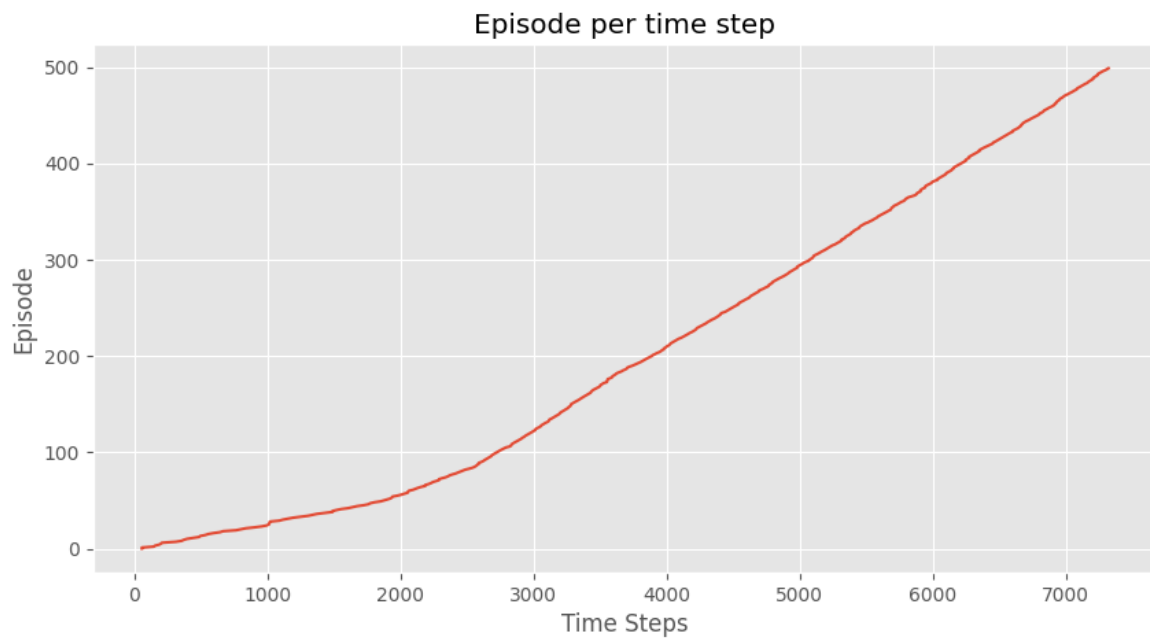
```
1 action = np.random.choice(np.arange(env.action_space.n), p=double_q_policy  
    (state))  
2 new_state, reward, done, _ = env.step(action)  
3 stats.episode_rewards[i_episode] += reward  
4 stats.episode_lengths[i_episode] = t
```

step 2 : TD Update:

```
1 ▼ if random.random() >= 0.5:
2     next_action = np.argmax(Q1[new_state])
3     Q1[state][action] += alpha * (reward + discount_factor * Q2[new_state]
4     [next_action] - Q1[state][action])
5 ▼ else:
6     next_action = np.argmax(Q2[new_state])
7     Q2[state][action] += alpha * (reward + discount_factor * Q1[new_state]
8     [next_action] - Q2[state][action])
9 ▼ if done:
10     break
11 state = new_state
```

- Corresponding Results:





## 2. 实现double-Q learning算法

Main implementation codes:

```

1 def make_double_epsilon_greedy_policy(Q1, Q2, epsilon, nA):
2     def policy_fn(observation):
3         A = np.ones(nA, dtype=float) * epsilon / nA
4         best_action = np.argmax(Q1[observation] + Q2[observation])
5         A[best_action] += (1.0 - epsilon)
6         return A
7     return policy_fn
8
9
10 def double_q_learning(env, num_episodes, discount_factor=1.0, alpha=0.5, epsilon=0.1):
11
12     # The final action-value function.
13     # A nested dictionary that maps state -> (action -> action-value).
14     Q1 = defaultdict(lambda: np.zeros(env.action_space.n))
15     Q2 = defaultdict(lambda: np.zeros(env.action_space.n))
16
17     # Keeps track of useful statistics
18     stats = plotting.EpisodeStats(
19         episode_lengths=np.zeros(num_episodes),
20         episode_rewards=np.zeros(num_episodes))
21
22     # The policy we're following
23     double_q_policy = make_double_epsilon_greedy_policy(Q1, Q2, epsilon, env.action_space.n)
24
25     for i_episode in range(num_episodes):
26         # Print out which episode we're on, useful for debugging.
27         if (i_episode + 1) % 100 == 0:
28             print("\rEpisode {}/{}.".format(i_episode + 1, num_episodes),
29 end="")
30             sys.stdout.flush()
31
32         # Reset the environment and pick the first action
33         state = env.reset()
34         for t in itertools.count():
35             #####Implement your code here#####
36             #####
37             # step 1 : Take a step
38             action = np.random.choice(np.arange(env.action_space.n), p=double_q_policy(state))
39             new_state, reward, done, _ = env.step(action)

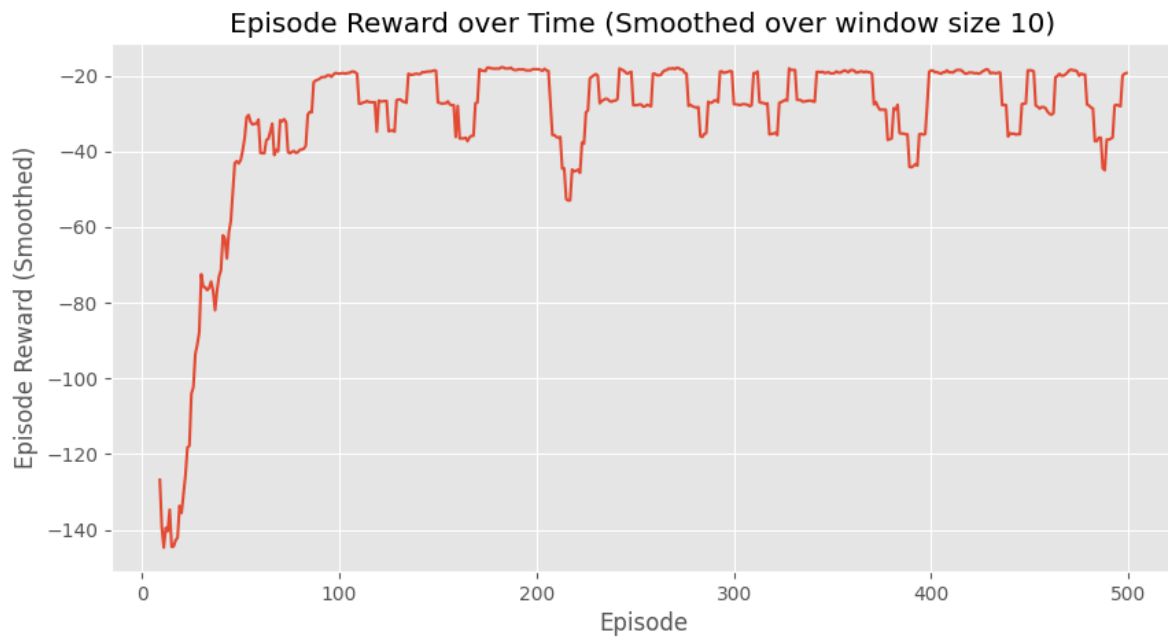
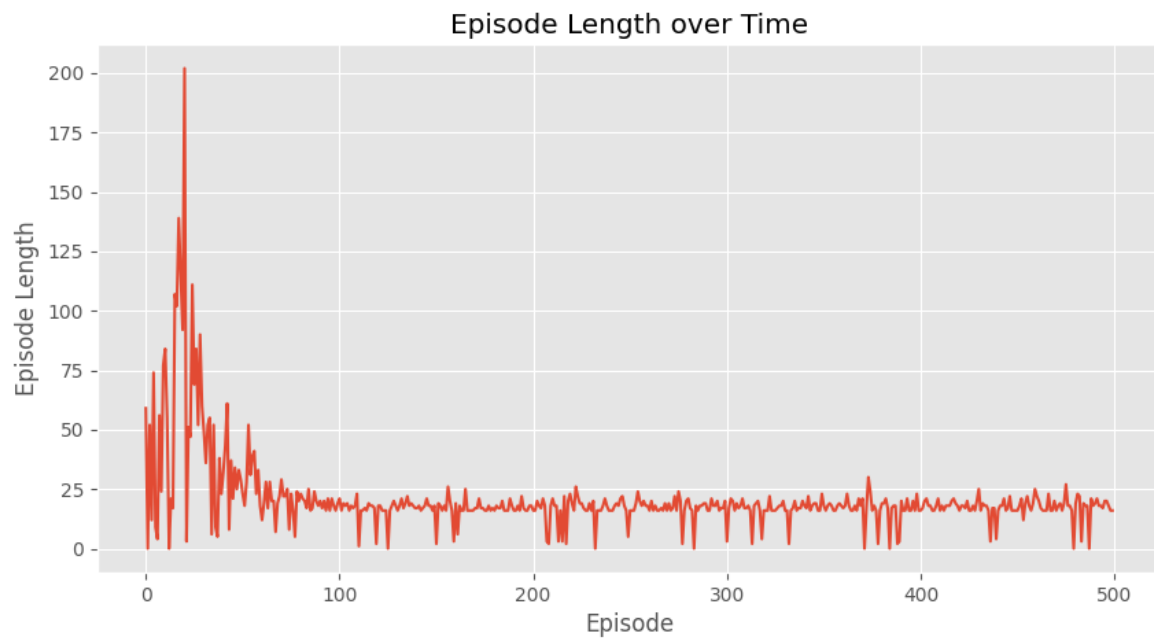
```

```

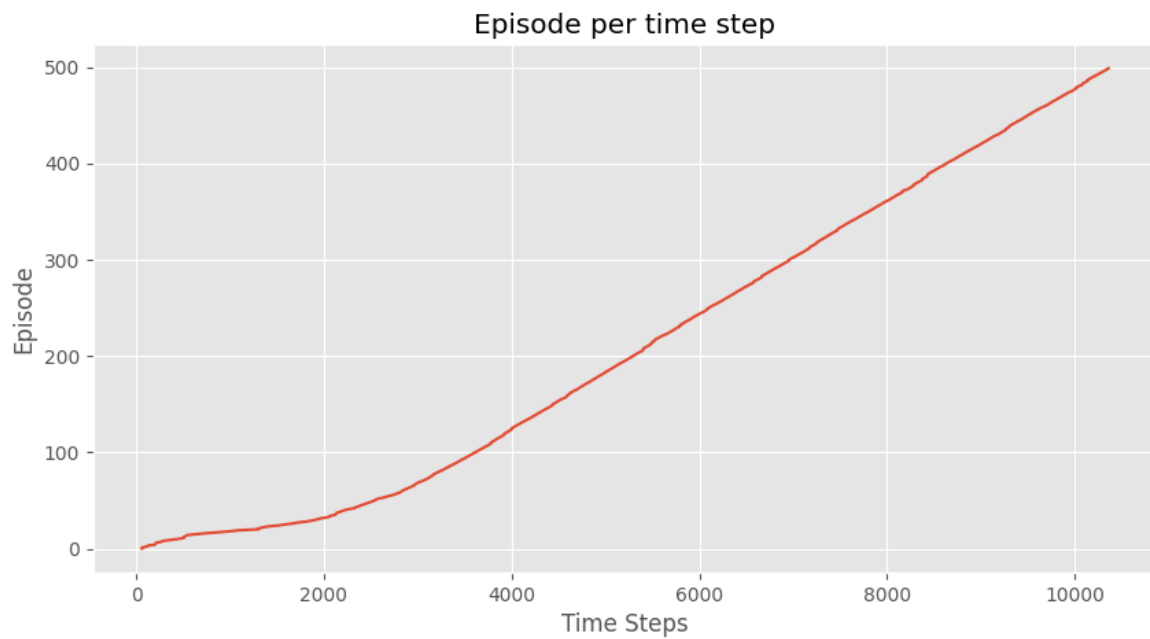
40
41         # Update statistics
42         stats.episode_rewards[i_episode] += reward
43         stats.episode_lengths[i_episode] = t
44
45     # step 2 : TD Update
46     if random.random() >= 0.5:
47         next_action = np.argmax(Q1[new_state])
48         Q1[state][action] += alpha * (reward + discount_factor * Q
49     2[new_state][next_action] - Q1[state][action])
50     else:
51         next_action = np.argmax(Q2[new_state])
52         Q2[state][action] += alpha * (reward + discount_factor * Q
53     1[new_state][next_action] - Q2[state][action])
54     if done:
55         break
56     state = new_state
57
58     #####Implement your code end#####
59     #####
60     return Q1, Q2, stats

```

- Corresponding Results







- Comparison between Q-learning & double-Q learning: Q-learning在每次选择动作时总是选择具有最大Q值的动作，但如果该动作是一个噪声动作，或者奖励是随机值时，可能会导致最大化偏差。这个问题的一个原因是在优化时使用了相同的样本。为了解决这个问题，Double Q-learning使用了两组独立的样本进行估计。根据结果分析，Double Q-learning的奖励曲线通常比Q-learning的曲线更加稳定。