

AI6101 Reinforcement Learning Assignment

Chen Lei

G2202273D

Nanyang Technological University

October 16, 2023

Abstract

This assignment implements a game agent to solve the CliffBoxPushing grid-world game, using Q-Learning algorithm. The Q-Learning algorithm adapt cosine annealing scheduler to decreased the ϵ parameter, which proves to be effective to get a relative good result.

Q-Learning

In this assignment, I choose Q-Learning algorithm to train the game agent. It is a model-free algorithm learning optimal policies based on reward observed when an action is taken at a specific state. The results(Q-values) are updated in a state-action table called Q-table. The Q-learning algorithm can be formulated as:

$$Q_{new}(S_t, A_t) = Q_{old}(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_{\alpha} Q_{old}(S_{t+1}, \alpha) - Q_{old}(S_t, A_t)) \quad (1)$$

where $Q_{new}(S_t, A_t)$ and $Q_{old}(S_t, A_t)$ refers to the new and old estimation when taking action A_t at state S_t , and α and γ is the learning rate and γ is the discount parameter.

Also, in order to encourage the agent to explore more available paths, the Epsilon Greedy Method is used to enable agent to take actions randomly. However, Q-value finally converges to a consistent value, so I uses cosine annealing scheduler, which can be formulated as:

$$\eta_t = \eta_{min}^i + \frac{1}{2}(\eta_{max}^i - \eta_{min}^i)(1 + \cos(\frac{T_{cur}}{T_i}\pi)) \quad (2)$$

where η_{max}^i is the largest ϵ and η_{min}^i is the minimum ϵ (0 in this scenario) and T_{cur} is the current episode during training.

Implementation and Training

In this assignment, the agent is trained for 10,000 episodes. I have tried to use scheduler on ϵ and learning rate α , which is shown in Figure1. It seems that using cosine annealing scheduler can not improve the rewards, even worse, the curve keeps vibrating at a

quite reward of $-1,602$ while the default Q-Learning agent gains the reward of -889 . However, the decay of ϵ seems to be helpful as the agent succeeds to converge after trained for nearly 6,000 episodes and gain 642 rewards at last. The final rewards are shown in table1.

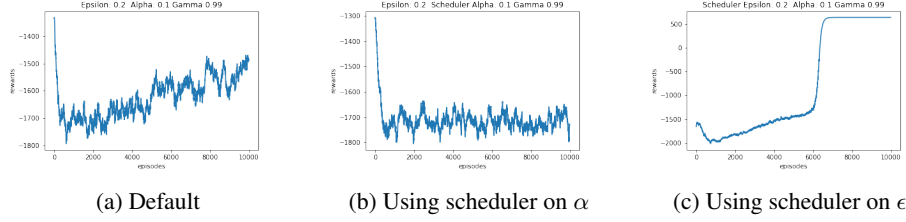


Figure 1: Episodes vs Rewards with three agents

Table 1: Rewards of three agents

Agent	Reward
Default	-889
Using α Scheduler	-1,602
Using ϵ Scheduler	642

After adapting cosine annealing scheduler on ϵ parameter during the training, I try to find a proper initial ϵ value. The experiments are conducted with 9 different ϵ value varying from 0.1 to 0.9 adding 0.1 each time. The episode vs reward curves of these experiments has little difference as figure3, which may means that a random initial ϵ value is suitable when using cosine annealing scheduler.

Result

After conducting the above experiments, I decide to train my agent for 10,000 episodes with initial $\alpha = 0.1$, $\gamma = 0.99$ and $\epsilon = 0.2$. I use V-table to show the best value of each state (here is the position of agent in this game) when applying the policy. The V-table is as the Figure2:

```
V table:
0 [-40.02, -40.31, -40.68, -40.94, -38.61, -44.45, 0.0, 0.0, -37.15, -31.7, -30.9, -29.78, -25.79, -24.35]
1 [-38.06, -35.75, -36.46, -44.64, -37.13, -43.3, 0.0, 0.0, -36.82, -29.51, -28.6, -27.73, -31.38, -23.94]
2 [-38.27, -36.83, -45.19, 0.0, -46.44, -44.03, 0.0, 0.0, -38.52, -30.51, -29.9, -41.43, 0.0, -28.37]
3 [-40.46, -35.84, -47.75, 0.0, -43.72, -43.6, 0.0, -44.51, -34.68, -34.11, -35.65, 0.0, 0.0, -30.22]
4 [-40.47, -37.27, -45.19, 0.0, -41.21, -35.8, -40.59, -33.72, -32.64, -31.97, -36.26, 0.0, 0.0, -28.71]
5 [-38.31, -38.35, -49.77, 0.0, -45.02, -38.48, -37.84, -36.82, -35.88, -32.53, -37.12, 0.0, 0.0, -27.63]
```

Figure 2: V-table

The agent can get the reward of 642 at last with the behavior history of [4, 1, 1, 1, 3, 1, 4, 4, 4, 4, 1, 4, 2, 2, 2, 3, 2, 4, 4, 4, 4, 4, 2, 4, 1, 1, 1, 3, 1, 4, 4, 4, 1, 4, 2, 2, 2] where 1, 2, 3, 4 refers to left, right, down and up. The policy is shown in the appendix result policy.

Appendix

Episode vs Rewards with different ϵ value

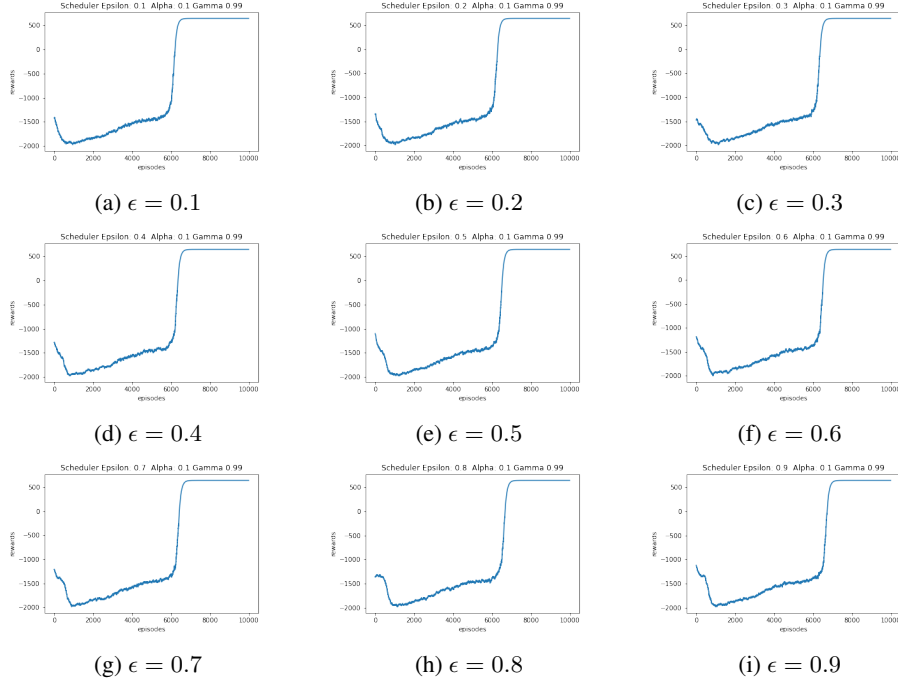


Figure 3: Episodes vs Rewards with different initial ϵ value

Result Policy

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]
[b'_' b'B' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G' ]
[b'A' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]
```

Learned Policy:

step: 1, state: (5, 0, 4, 1), actions: 4, reward: -14

Action: 4

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]
[b'_' b'B' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G' ]]
```



```

[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'x' b'A']
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'B']
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G']
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]
step: 37, state: (2, 13, 3, 13), actions: 2, reward: 998
Action: 2
[[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_' b'_' ]
[b'_' b'_' b'_' b'x' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'A']
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'B']
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]

```