

CLSQL: Continuous Local Search Q-Learning

Improved Q-Learning Algorithm Based on Continuous Local Search Policy for Mobile Robot Path Planning

Abdallah Amr Mostafa Osama Tarek Shohdy Zeiad
El-shabrawy

Faculty of Engineering
Egypt Japan University for Science and Technology
MTE311 Seminar on MTE
Dr. Haitham El-Hussieny

February 27, 2023

Table of Contents

1 Continuous Local Search Q-Learning Algorithm

- Environment Prior Knowledge
- Local Search Policy
- Dynamic ϵ -greedy
- Optimal Steps
- Search and Iteration Process

2 Implementation

Motivation

Problems in the QL algorithm:

- slow convergence speed
- large number of iterations

Improved Q-learning algorithm based on a continuous local search policy, CLSQL.

- 1 A prior knowledge to initialize the Q-table.
- 2 Adding the local search policy.
- 3 Adjusting ϵ -greedy policy.

Environment Prior Knowledge

- Q-tables in the initial state are equal to 0.
- Distance function is used to determine the prior knowledge of the environment.

$$D_s = \sqrt{(x_E - x)^2 + (y_E - y)^2}$$

$$Q(s, a) = \begin{cases} \frac{1}{D_s} & D_s > 0 \\ 0 & D_s = 0 \end{cases}$$

Environment Prior Knowledge

Algorithm 1: Prior Knowledge

- 1 Initialize Agent State $s(x, y)$, Destination State $s_E(x_E, y_E)$, Obstacle State s_{Obs}
- 2 Repeat
- 3 if $s == s_{Obs}$ then
- 4 $D_s = 0$
- 5 else
- 6 $D_s = \sqrt{(x_E - x)^2 + (y_E - y)^2}$
- 7 Until All s traversal completed

Local Search Policy

- As the size of the environment increases, dimensions increase.
- Large number of iterations.

Local search policy can simplify the complex environment by:

- 1 By setting a Local Environment based on s_{start} or s_I .
- 2 Centering of s_{start} or s_I and obtain the local environment grid.
- 3 s_I are determined based on prior knowledge.

Local Search Policy

Algorithm 2: Local Environment

```
1 Initialize  $LG, LS, s_{start(x,y)}, s_l(x,y)$ 
2 for  $i = \lceil -\frac{LS}{2} \rceil, \dots, 0, \dots, \lfloor \frac{LS}{2} \rfloor$  do
3    $j = \lceil -\frac{LS}{2} \rceil, \dots, 0, \dots, \lfloor \frac{LS}{2} \rfloor$  do
4     if  $i + x \geq 0, \&i + x < G_S, \&j + y \geq 0, \&j + y < G_S$  then
5        $LG_{i+x,j+y} \leftarrow GG_{i+x,j+y}$ 
6   end for
7 end for
```

Local Search Policy

Search effectiveness were increased by determining s_l in the local environment.

We use the Priority Queue (PriQ) to determine the priority of s_l

$$PriQ = D_s + \frac{GS}{|r_{obs}|} P_s$$

GS is the Global Environment Size

P_s is the sum of reward in eight neighborhoods of the state s ,

Local Search Policy

Algorithm 3: Intermediate Points

```
1 Initialize  $s_E, s_I$ 
2 if  $s_E \in LG$  then
3    $s_I = s_E$ 
4 else
5    $PriQ = D_s + \frac{GS}{|r_{obs}|} P_s$ 
6    $s_I = argmin(PriQ)$ 
7 Return  $s_I$ 
```

Local Search Policy

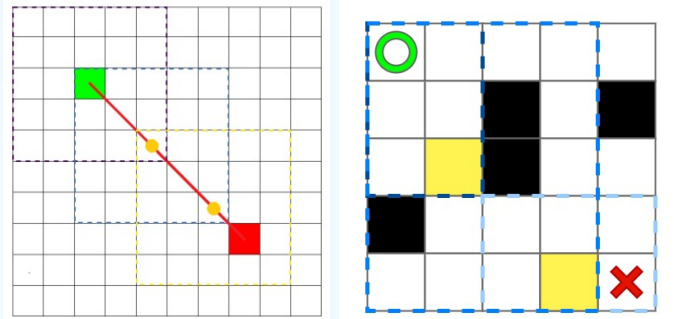


Figure: Grid maps with intermediate point.

Dynamic ϵ -greedy

Adjusted ϵ -greedy policy:

- 1 ϵ is initialized at the beginning of each local environment search.
- 2 ϵ increases when the agent collides with obstacles.
- 3 ϵ decreases when the result is the non-optimal path.

Optimal Steps

when number of learning steps at this stage is far greater than the local environment size ($steps > LS^2$), optimal solution cannot be obtained.

Present s_l in the local environment is deleted, and then another new s_l is determined again.

Global optimization can be achieved only when local optimization is achieved.

The optimal steps correspond to the diagonal distance between two points p_1 and p_2 .

$$D = |x_2 - x_1| + |y_2 - y_1| + (\sqrt{2} - 2) * \min(|x_2 - x_1|, |y_2 - y_1|)$$

Search and Iteration Process

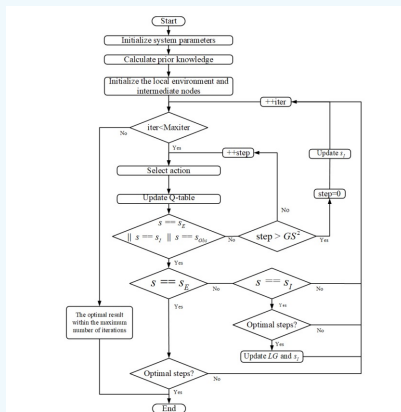


Figure: The flow of the CLSQL algorithm

Table of Contents

- 1 Continuous Local Search Q-Learning Algorithm
 - Environment Prior Knowledge
 - Local Search Policy
 - Dynamic ϵ -greedy
 - Optimal Steps
 - Search and Iteration Process
- 2 Implementation

Implementation

```
1 import numpy as np
2 import math
3 import gym
4 import random
5 from time import sleep as sleep
6
7 env = gym.make("FrozenLake8x8-v1", is_slippery=False)
8
9 n_actions = env.action_space.n
10 n_states = env.observation_space.n
11
12 n_episodes = 40
13 max_steps = 275
14
15 gamma = 0.99
16 learning_rate = 0.86
17
18 epsilon = 0.99
19 max_epsilon = 1.0
20 min_epsilon = 0.01
21 decay = 0.1
22
23 rewards = []
24
25 q_table = np.zeros((n_states, n_actions))
26
27
```

Implementation

```
1 for s in range(n_states):
2     if s in {35, 41, 42, 46, 49, 19, 52, 54, 59, 29, 63}:
3         q_table[s,a] = 0
4         continue
5     for a in range(n_actions):
6         if a == 0:
7             st = s - 1
8             if st < 0:
9                 continue
10        if a == 1:
11            st = s + 8
12            if st < 0:
13                continue
14        if a == 2:
15            st = s + 1
16            if st < 0:
17                continue
18        if a == 3:
19            st = s - 8
20            if st < 0:
21                continue
22        x = int(st / 8)
23        y = st % 8
24        if x == 7 and y == 7:
25            q_table[s,a] = 10
26            continue
27        q_table[s,a] = 1 / math.sqrt((7-x)**2+(7-y)**2)
```


Implementation

```
1 for episode in range(n_episodes+1):
2     state = env.reset()
3     step = 0
4     done = False
5     total_rewards = 0
6     for step in range(max_steps):
7
8         e_e_tradeoff = random.uniform(0,1)
9         if(e_e_tradeoff >= epsilon):
10             action = np.argmax(q_table[state,:])
11         else:
12             action = env.action_space.sample()
13         new_state, reward, done, info = env.step(action)
14         if reward == 0 and done == True:
15             q_table[state,action] = q_table[state,action] + learning_rate * ((
16 reward-10) + gamma * np.max(q_table[new_state,:]) - q_table[state,action])
17         elif reward == 1 and done == True:
18             q_table[state,action] = q_table[state,action] + learning_rate * ((
19 reward+9) + gamma * np.max(q_table[new_state,:]) - q_table[state,action])
20         else:
21             q_table[state,action] = q_table[state,action] + learning_rate * ((
22 reward) + gamma * np.max(q_table[new_state,:]) - q_table[state,action])
23
24         total_rewards += reward
25         state = new_state
26         if(done == True):
27             break
```