

## Assignment (강화학습 - 김기응 교수님 연구실)

담당조교: 이종민 (jmlee@ai.kaist.ac.kr)

### -과제설명

실습시간에 구현해본 Q-Learning을 확장한 **Watkin's  $Q(\lambda)$ 알고리즘**을 구현해보는 과제입니다. 매 타임스텝마다 획득한 보상 정보를 토대로 Q 테이블을 업데이트 할 때, Q-Learning은 딱 한 개의  $Q(s,a)$ 값만 업데이트 시키는 점을 기억하시나요? 하지만 이 경우 maze 도메인에서처럼 Goal 상태에서 받는 보상 정보가 초기 상태의 Q값에 전달되어 반영될 때까지는 많은 시간을 필요로 하게 됩니다.

$Q(\lambda)$  알고리즘은 에이전트의 이동해온 경로, trace를 정보를 이용해 Q 테이블의 여러  $(s,a)$  값을 한 번에 업데이트 시킵니다. 이러면 Goal 지점에서 받은 보상 신호가 기존 Q-Learning에 비해 훨씬 빠르게 다른  $(s,a)$ 들로 퍼져나가겠죠?

다음 페이지 Watkin's  $Q(\lambda)$  알고리즘의 Pseudo-code를 참고하여 구현해주시면 됩니다. (Q-Learning의 Pseudo-code는 참고 및 비교용으로 함께 첨부드렸습니다.)

### -제출물

완성된 코드(`assignment_rl.py`)와 보고서 (docx나 pdf로 형식)를 제출해 주시면 됩니다

### -채점

채점은 코드 5점, 보고서 5점, 총 10점 만점으로 채점됩니다.

보고서에는

- Maze 환경에서 (다른 도메인의 결과도 포함하셔도 좋습니다) Q-Learning과  $Q(\lambda)$ 을 각각 어떤 결과의 차이가 있는지
- 다양한 파라미터 (예: epsilon, alpha 등)를 바꿔가며 실행했을 때, 결과의 차이가 있는지. 있다면 어떻게 달라지는지

를 서술해주시면 됩니다.

---

**Algorithm 1** Q-Learning

---

```
1:  $Q(s, a) \leftarrow 0$  for all  $s, a$ 
2: for episode = 0, 1, ... do
3:    $s \leftarrow \text{env.reset}()$ 
4:    $a \leftarrow \text{EPSILONGREEDY}(Q(s, \cdot))$ 
5:   for  $t = 0, 1, \dots$  do
6:      $(s', r, \text{done}) \leftarrow \text{env.step}(a)$ 
7:      $a' \leftarrow \text{EPSILONGREEDY}(Q(s', \cdot))$ 
8:      $a^* \leftarrow \arg \max_b Q(s', b)$ 
9:     if done then
10:       $\delta \leftarrow r - Q(s, a)$ 
11:     else
12:       $\delta \leftarrow r + \gamma Q(s', a^*) - Q(s, a)$ 
13:     end if
14:      $Q(s, a) \leftarrow Q(s, a) + \alpha \delta$ 
15:     if done then
16:       break
17:     end if
18:      $s \leftarrow s'$  and  $a \leftarrow a'$ 
19:   end for
20: end for
```

---

---

**Algorithm 2** Watkin's  $Q(\lambda)$ 

---

```
1:  $Q(s, a) \leftarrow 0$  for all  $s, a$ 
2: for episode = 0, 1, ... do
3:    $e(s, a) \leftarrow 0$  for all  $s, a$ 
4:    $s \leftarrow \text{env.reset}()$ 
5:    $a \leftarrow \text{EPSILONGREEDY}(Q(s, \cdot))$ 
6:   for  $t = 0, 1, \dots$  do
7:      $(s', r, \text{done}) \leftarrow \text{env.step}(a)$ 
8:      $a' \leftarrow \text{EPSILONGREEDY}(Q(s', \cdot))$ 
9:      $a^* \leftarrow \arg \max_b Q(s', b)$ 
10:    if done then
11:       $\delta \leftarrow r - Q(s, a)$ 
12:    else
13:       $\delta \leftarrow r + \gamma Q(s', a^*) - Q(s, a)$ 
14:    end if
15:     $e(s, a) \leftarrow e(s, a) + 1$ 
16:    for all  $s, a$  do
17:       $Q(s, a) \leftarrow Q(s, a) + \alpha \delta e(s, a)$ 
18:      if  $a' = a^*$  then
19:         $e(s, a) \leftarrow \gamma \lambda e(s, a)$ 
20:      else
21:         $e(s, a) \leftarrow 0$ 
22:      end if
23:    end for
24:    if done then
25:      break
26:    end if
27:     $s \leftarrow s'$  and  $a \leftarrow a'$ 
28:  end for
29: end for
```

---

- $Q(\lambda)$  알고리즘에는 trace를 얼마나 빨리 감쇠 시킬지 결정하는  $\lambda$  파라미터가 하나 더 존재합니다. ( $\lambda = 0$ ) 이면  $Q(\lambda)$ 와 Q-Learning 알고리즘은 동일해집니다.