

## Assignment (강화학습 - 김기응 교수님 연구실)

담당교교: 이병준, 윤든솔 (bjlee@ai.kaist.ac.kr, dsyoon@ai.kaist.ac.kr)

### -과제설명

실습시간에 구현해본 SARSA를 확장한 **SARSA( $\lambda$ )알고리즘을 구현**해보는 과제입니다. 매 타임스텝마다 획득한 보상 정보를 토대로 Q 테이블을 업데이트 할 때, SARSA은 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.

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

다음 페이지 SARSA( $\lambda$ ) 알고리즘의 Pseudo-code를 참고하여 구현해주시면 됩니다. (SARSA의 Pseudo-code는 참고 및 비교용으로 함께 첨부드렸습니다.)

### -필요환경

numpy, gym, pygame

### -제출물

완성된 코드(assignment.py)를 제출해 주시면 됩니다

### -채점

채점은 코드 10점 만점으로 채점됩니다.

---

**Algorithm 1** SARSA

---

```
1:  $Q(\tilde{s}, \tilde{a}) \leftarrow 0$  for all  $(\tilde{s}, \tilde{a}) \in \mathcal{S} \times \mathcal{A}$ 
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:     if done then
9:        $\delta \leftarrow r - Q(s, a)$ 
10:    else
11:       $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$ 
12:    end if
13:     $Q(s, a) \leftarrow Q(s, a) + \alpha \delta$ 
14:    if done then
15:      break
16:    end if
17:     $s \leftarrow s'$  and  $a \leftarrow a'$ 
18:  end for
19: end for
```

---

---

**Algorithm 2** SARSA( $\lambda$ )

---

```
1:  $Q(\tilde{s}, \tilde{a}) \leftarrow 0$  for all  $(\tilde{s}, \tilde{a}) \in \mathcal{S} \times \mathcal{A}$ 
2: for episode = 0, 1, ... do
3:    $e(\tilde{s}, \tilde{a}) \leftarrow 0$  for all  $(\tilde{s}, \tilde{a}) \in \mathcal{S} \times \mathcal{A}$ 
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:     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:     $e(s, a) \leftarrow e(s, a) + 1$ 
15:    for all  $(\tilde{s}, \tilde{a}) \in \mathcal{S} \times \mathcal{A}$  do
16:       $Q(\tilde{s}, \tilde{a}) \leftarrow Q(\tilde{s}, \tilde{a}) + \alpha \delta e(\tilde{s}, \tilde{a})$ 
17:       $e(\tilde{s}, \tilde{a}) \leftarrow \gamma \lambda e(\tilde{s}, \tilde{a})$ 
18:    end for
19:    if done then
20:      break
21:    end if
22:     $s \leftarrow s'$  and  $a \leftarrow a'$ 
23:  end for
24: end for
```

---

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