# SARSA vs Q-learning

SARSA Algorithm

$\epsilon$-greedy policy improvement done for TD methods

1: Set initial $\epsilon$-greedy policy $\pi$, $t = 0$, initial state $s_t = s_0$

2: Take $a_t \sim \pi(s_t)$ // Sample action from policy

3: Observe $(r_t, s_{t+1})$

4: **loop**

5:     Take action $a_{t+1} \sim \pi(s_{t+1})$

6:     Observe $(r_{t+1}, s_{t+2})$

7:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$

8:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random

E-greedy policy update?

9:     $t = t + 1$

10: **end loop**

Loop Q value update for every time-step : based on TD methods

Sutton
and
Barto

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

Action chosen by e-greedy

Action taken for next state already decided

1: Set initial $\epsilon$-greedy policy $\pi$, $t = 0$, initial state $s_t = s_0$
2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
3: Observe $(r_t, s_{t+1})$
4: **loop**
5:     Take action $a_{t+1} \sim \pi(s_{t+1})$
6:     Observe $(r_{t+1}, s_{t+2})$
7:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
8:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
9:     $t = t + 1$
10: **end loop**

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

1: Set initial $\epsilon$-greedy policy $\pi$, $t = 0$, initial state $s_t = s_0$
2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
3: Observe $(r_t, s_{t+1})$
4: **loop**
5:    Take action $a_{t+1} \sim \pi(s_{t+1})$
6:    Observe $(r_{t+1}, s_{t+2})$
7:    $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
8:    $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
9:    $t = t + 1$
10: **end loop**

E-greedy policy update?

Loop Q value update for every time-step : based on TD methods

---

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

Action chosen by e-greedy

Action taken for next state already decided

Difference?

Sutton
and
Barto

Seems to indicate that actions are taken greedily

No indication of how next action is chosen

But can assume greedy action choosing

1: Set initial $\epsilon$-greedy policy $\pi$, $t = 0$, initial state $s_t = s_0$
2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
3: Observe $(r_t, s_{t+1})$
4: **loop**
5:     Take action $a_{t+1} \sim \pi(s_{t+1})$
6:     Observe $(r_{t+1}, s_{t+2})$
7:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
8:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
9:     $t = t + 1$
10: **end loop**

E-greedy policy update?

Loop Q value update for every time-step based on TD methods

But this would mean q-value has been calculated for all actions at state st

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

Action chosen by e-greedy

Action taken for next state already decided

lecture 4

stop

Simple difference between SARSA and Q-learning

SARSA = on-policy
Q-learning = off-policy

# SARSA = on-policy

## on-policy = on e-greedy policy

## on-policy = learn to estimate and evaluate policy from experience obtained from following that policy (purely from experience)

Sutton
and
Barto

1: Set initial $\epsilon$-greedy policy $\pi$, $t = 0$, initial state $s_t = s_0$
2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
3: Observe $(r_t, s_{t+1})$
4: **loop**
5:     Take action $a_{t+1} \sim \pi(s_{t+1})$
6:     Observe $(r_{t+1}, s_{t+2})$
7:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
8:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
9:     $t = t + 1$
10: **end loop**

E-greedy policy update?

Loop Q value update for every time-step : based on TD methods

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

Action chosen by e-greedy

Action taken for next state already decided

**SARSA = on-policy**

**on-policy = on e-greedy policy**

**on-policy = learn to estimate and evaluate policy from experience obtained from following that policy (purely from experience)**

1: Set initial $\epsilon$-greedy policy $\pi$, $t = 0$, initial state $s_t = s_0$
2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
3: Observe $(r_t, s_{t+1})$
4: **loop**
5:     Take action $a_{t+1} \sim \pi(s_{t+1})$
6:     Observe $(r_{t+1}, s_{t+2})$
7:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
8:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
9:     $t = t + 1$
10: **end loop**

E-greedy policy update?

Loop Q value update for every time-step : based on TD methods

**Seems to be experience following a different policy (trying out all other actions)**

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R$, $S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

Action chosen by e-greedy

Action taken for next state already decided

**The policy on the lefthandside seems stochastic?**

Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)

8:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random

**Action already pre-determined**

**When choosing action from this policy**

**If pi takes a single action = either argmax or random by epsilon**

**Then**

**pi = Deterministic policy**

lecture 4

Conditions of SARSA:

GLIE
Robbin-Munro

lecture 4

Conditions of SARSA:

GLIE
Robbin-Munro

## Definition of GLIE

- All state-action pairs are visited an infinite number of times

$$\lim_{i\to\infty} N_i(s, a) \to \infty$$

- Behavior policy converges to greedy policy
$\lim_{i\to\infty} \pi(a|s) \to \arg\max_a Q(s, a)$ with probability 1

Basically being able to
explore all actions

= use method of e-
greedy where e starts
big and then diminishes

All s, a are visited
infinitely often

lecture 4

Conditions of SARSA:

GLIE
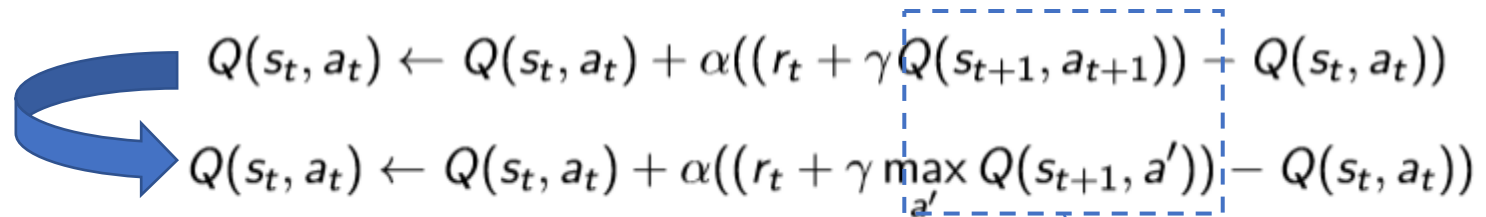Robbin-Munro

② The step-sizes $\alpha_t$ satisfy the Robbins-Munro sequence such that

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Learning rate not too
big..?

lecture 4

Conditions of SARSA:

GLIE
Robbin-Munro

2 The step-sizes $\alpha_t$ satisfy the Robbins-Munro sequence such that

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Learning rate not too
big..?

Alpha = 0 : no learning

Alpha = 1 : new experience
complete overtakes

lecture 4

SARSA → Q-Learning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t))$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

Bootstrapping

1: Initialize $Q(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$

2: Set $\pi_b$ to be $\epsilon$-greedy w.r.t. $Q$

3: **loop**

4:     Take $a_t \sim \pi_b(s_t)$ // Sample action from policy

5:     Observe $(r_t, s_{t+1})$

6:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \arg\max_a Q(s_{t_1}, a) - Q(s_t, a_t))$

7:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random

8:     $t = t + 1$

9: **end loop**

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\textit{terminal}, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        Take action $A$, observe $R$, $S'$
        $Q(S, A) \leftarrow Q(S, A) + \alpha\big[R + \gamma \max_a Q(S', a) - Q(S, A)\big]$
        $S \leftarrow S'$
    until $S$ is terminal

lecture 4

**Q-learning = off policy**

**Off-policy = estimate and evaluate policy from following a different policy**

---

1: Initialize $Q(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
2: Set $\pi_b$ to be $\epsilon$-greedy w.r.t. $Q$
3: **loop**
4:     Take $a_t \sim \pi_b(s_t)$ // Sample action from policy
5:     Observe $(r_t, s_{t+1})$
6:     $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \arg\max_a Q(s_{t_1}, a) - Q(s_t, a_t))$
7:     $\pi(s_t) = \arg\max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
8:     $t = t + 1$
9: **end loop**

---

**Not necessarily the same policy**

lecture 4

**Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        Take action $A$, observe $R, S'$
        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma \max_a Q(S', a) - Q(S, A) \big]$
        $S \leftarrow S'$
    until $S$ is terminal

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma Q(S', A') - Q(S, A) \big]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

**Not necessarily the same policy**

**Difference = random by epsilon?**

# lecture 4



**Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Algorithm parameters: step size $\alpha \in (0,1]$, small $\varepsilon > 0$
Initialize $Q(s,a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        Take action $A$, observe $R, S'$
        $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma \max_a Q(S',a) - Q(S,A)]$
        $S \leftarrow S'$
    until $S$ is terminal

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0,1]$, small $\varepsilon > 0$
Initialize $Q(s,a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma Q(S',A') - Q(S,A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

**Not necessarily the same policy**

**Difference =**

**No matter what next action chosen for evaluating current state,**

**Next action will be decided with the updated policy (new and different from that used to evaluate)**

**When evaluated with next action chosen, the chosen next action will be used on next step**

Q-learning tends to have better results

For most cases...
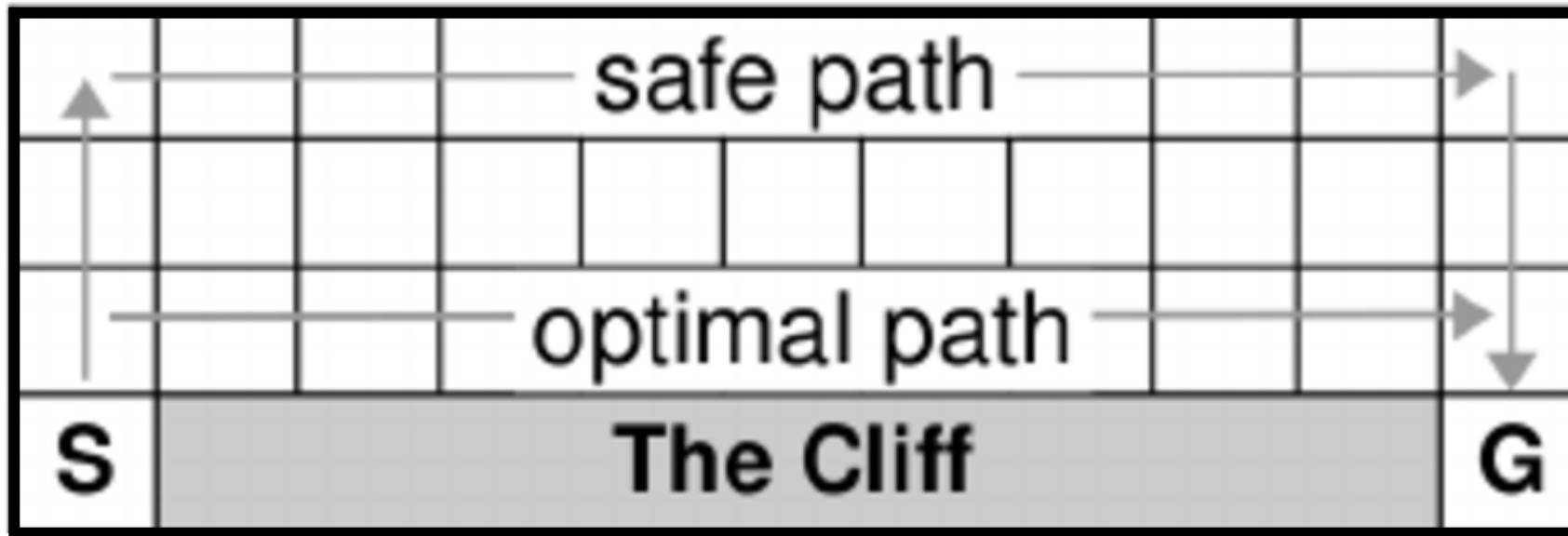
But for some cases...

Why SARSA rather than Q-learning?

Sutton and Barto example of cliff-walking

SARSA may do better on early stages

SARSA is realistic
Q-learning is optimistic

Case where there could be a lot of negative outcomes from optimism of using max

Optimism would mean being only interested in the best outcome

Optimism on Cliff for example:

Near cliff is the shortest route to finish = highest possible reward & but forgets that it also has high chance of falling

# SARSA Code breakdown

```
 9   import gym
10   import gym_gridworlds          ← Import package that makes &
11   import numpy as np               registers Cliff environment
12   import matplotlib.pyplot as plt
13   from gym.envs.registration import register
14   import random as pr
15
16
17   env = gym.make('Cliff-v0')     ← Create element
```

# SARSA Code breakdown

Q value (state-action value) is defined by value that comes from action taken at current state

= preparing Table for all possible states and actions

```
19    # Q-table
20    # note that:
21    # action = [up, right, down, left]
22    Q = np.zeros([env.height , env.width, env.action_space.n])
```

# SARSA Code breakdown

```
19   # Q-table
20   # note that:
21   # action = [up, right, down, left]
22   Q = np.zeros([env.height , env.width, env.action_space.n])
```

Current state can be defined by current position (x, y)

## SARSA Code breakdown

```
24    # define hyperparameters (epsilon defined as function later)
25    dis = .91
26    num_episodes = 1000
27    lr = 0.01
28    success = 0
29
30    # create lists to contain total rewards and steps per episode
31    rList = []
```

= to track how total reward changes

# SARSA Code breakdown

```python
def reset(self):
    self.S = (3, 0)
    return self.S
```

```python
32    for i in range(num_episodes):
33        # Reset environment and get first new ob
34        (y, x) = env.reset()
35        # state = y , x
36        rAll = 0
37        done = False
38
39        e = 1. / ((i // 100) + 1)
40        # e-greedy where epsilon starts big
41        # and becomes small
42        # = more exploration in the beginning and less later on
43        # with larger num_episodes, this might come into better effect
44
45        # note that this is a method to satisfy condition of GLIE
```

| R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 |
| R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 | R= -1 |
| S | The Cliff (R = -100) | | | | | | | | G |

# SARSA Code breakdown

```
32    for i in range(num_episodes):
33        # Reset environment and get first new observation
34        (y, x) = env.reset()
35        # state = y , x
36        rAll = 0
37        done = False
38
39        e = 1. / ((i // 100) + 1)
40        # e-greedy where epsilon starts big
41        # and becomes small
42        # = more exploration in the beginning and less later on
43        # with larger num_episodes, this might come into better effect
44
45        # note that this is a method to satisfy condition of GLIE
```

Done = reached destination

When done = True, episode is terminated

As #episode trained increases,

Epsilon decreases

# SARSA Code breakdown

```
47        # e-greedy
48        if np.random.rand(1) < e:
49            action = env.action_space.sample()
50        else:
51            action = np.argmax(Q[y, x, :])
```

# SARSA Code breakdown

```python
53        # The Q-Table learning algorithm
54    while not done:
55        # Get new state and reward from environment
56        (new_y, new_x), reward, done, _ = env.step(action)
57        # new_state = new_x , new_y
58        # Choose the next action by e greedy
59        if np.random.rand(1) < e:
60            next_action = env.action_space.sample()
61        else:
62            next_action = np.argmax(Q[new_y, new_x, :])
```

# SARSA Code breakdown

```python
53        # The Q-Table learning algorithm
54        while not done:
55            # Get new state and reward from environment
56            (new_y, new_x), reward, done, _ = env.step(action)
```

```python
24        def step(self, action):
25            x, y = self.moves[action]
26            self.S = self.S[0] + x, self.S[1] + y
27
28            self.S = max(0, self.S[0]), max(0, self.S[1])
29            self.S = (min(self.S[0], self.height - 1),
30                      min(self.S[1], self.width - 1))
31
32            if self.S == (self.height - 1, self.width - 1):
33                return self.S, -1, True, {}
34            elif self.S[1] != 0 and self.S[0] == self.height - 1:
35                # the cliff
36                return self.S, -100, True, {}
37            return self.S, -1, False, {}
```

```python
14        self.moves = {
15            0: (-1, 0),   # up
16            1: (0, 1),    # right
17            2: (1, 0),    # down
18            3: (0, -1),   # left
19        }
```

# SARSA Code breakdown

```
64          # Update Q-Table with new knowledge using learning rate
65          Q[y, x, action] = (Q[y, x, action] + lr * (reward + dis * Q[new_y, new_x, next_action] - Q[y, x, action]))
66
67          rAll += reward
68          y, x = new_y, new_x
69
70          # difference with Q-learning
71          # action is already chosen for updating current state-action value
72          action = next_action
73          # therefore the loop continues by choosing only the next action
74
75      rList.append(rAll)
```

Evaluate using
next_action determined
by the current policy

# SARSA Code breakdown

```
64          # Update Q-Table with new knowledge using learning rate
65          Q[y, x, action] = (Q[y, x, action] + lr * (reward + dis * Q[new_y, new_x, next_action] - Q[y, x, action]))
66
67          rAll += reward
68          y, x = new_y, new_x
69
70          # difference with Q-learning
71          # action is already chosen for updating current state-action value
72          action = next_action
73          # therefore the loop continues by choosing only the next action
74
75      rList.append(rAll)
```

Then use the next_action for next_state in next timestep

# SARSA Code breakdown

```
64          # Update Q-Table with new knowledge using learning rate
65          Q[y, x, action] = (Q[y, x, action] + lr * (reward + dis * Q[new_y, new_x, next_action] - Q[y, x, action]))
66
67          rAll += reward
68          y, x = new_y, new_x
69
70          # difference with Q-learning
71          # action is already chosen for updating current state-action value
72          action = next_action
73          # therefore the loop continues by choosing only the next action
74
75      rList.append(rAll)
```

Loop until termination

# SARSA Code breakdown

```python
77    print("Success rate: " + str(sum(rList) / num_episodes))
78    print("Final Q-Table Values")
79    print(Q)
80    print("Q Values near starting position")
81    print("START\n" + str(Q[3,0:1]))
82    print("Above START\n" + str(Q[2,0:2]))
83    print("Above above START\n" + str(Q[1,0:2]))
84    plt.bar(range(len(rList)), rList, color="blue")
85    plt.show()
```

# Q-learning Code breakdown

```
31      # The Q-Table learning algorithm
32      while not done:
33
34          if np.random.rand(1) < e:
35              action = env.action_space.sample()      Second difference
36          else:
37              action = np.argmax(Q[y, x, :])
38          # Get new state and reward from environment
39          (new_y, new_x), reward, done, _ = env.step(action)
40          # new_state = new_x , new_y
41
42          # Update Q-Table with new knowledge using learning rate
                                                          Main difference
43          Q[y, x, action] = (Q[y, x, action] + lr * (reward + dis * np.max(Q[new_y, new_x, :]) - Q[y, x, action]))
44
45          rAll += reward
46          y, x = new_y, new_x
47
48      rList.append(rAll)
```

# Q-learning tends to have better results

SARSA

```
laphisboy@laphisboy-900X3K:~/RL_fall/fall_week_1$ python3 just_lake_SARSA.py
Success rate: 0.7905
Final Q-Table Values
[[6.70736658e-02 4.77326352e-01 5.80999368e-03 8.48964303e-02]
 [1.05166556e-01 0.00000000e+00 1.61637981e-07 3.68353915e-04]
 [1.03474654e-02 1.90521207e-07 1.18426275e-11 1.21421652e-06]
 [5.78286553e-07 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [7.22760373e-02 5.54426264e-01 0.00000000e+00 5.51919838e-02]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 7.07570326e-04]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [1.24071121e-01 0.00000000e+00 6.44242082e-01 9.22765749e-02]
 [9.58373591e-02 5.85828924e-02 7.64135946e-01 0.00000000e+00]
 [1.57588989e-01 8.84203921e-01 0.00000000e+00 3.71905143e-05]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 6.08529223e-03 4.09755896e-01 5.92824654e-03]
 [5.26202400e-02 2.62430749e-01 9.99999874e-01 1.93303018e-01]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]]
```
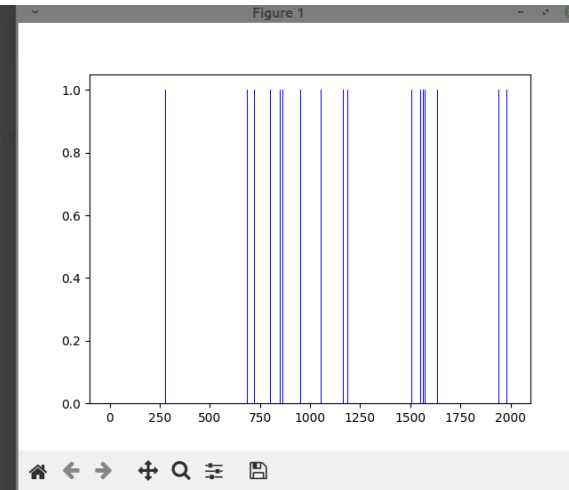
Q-Learning

```
laphisboy@laphisboy-900X3K:~/RL_fall/fall_week_1$ python3 just_lake_q.py
Score over time: 0.81
Final Q-Table Values
[[0.531441   0.59049    0.59049    0.531441  ]
 [0.531441   0.         0.6561     0.59048989]
 [0.59036441 0.729      0.4782969  0.64894266]
 [0.649539   0.         0.         0.        ]
 [0.59049    0.6561     0.         0.531441  ]
 [0.         0.         0.         0.        ]
 [0.         0.81       0.         0.64829366]
 [0.         0.         0.         0.        ]
 [0.6561     0.         0.729      0.59049   ]
 [0.6561     0.81       0.81       0.        ]
 [0.72899927 0.9        0.         0.72826444]
 [0.         0.         0.         0.        ]
 [0.         0.         0.         0.        ]
 [0.         0.81       0.9        0.729     ]
 [0.81       0.9        1.         0.81      ]
 [0.         0.         0.         0.        ]]
```

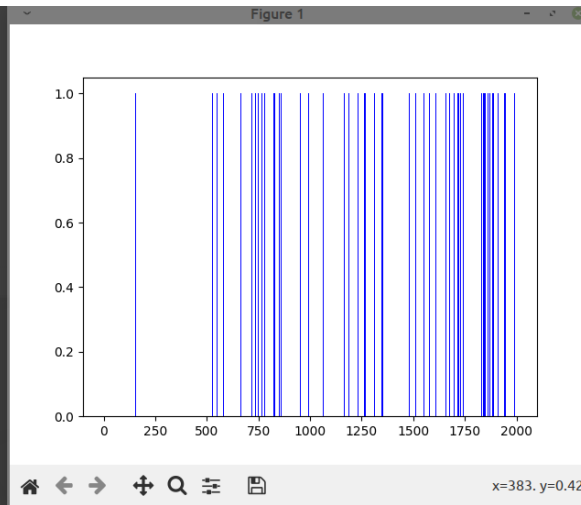# Q-learning tends to have better results



SARSA

```
laphisboy@laphisboy-900X3K:~/RL_fall/fall_week_1$ python3 frozen_lake_SARSA.py
Success rate: 0.0545
Final Q-Table Values
[[1.55087990e-03 6.78733437e-03 1.30792462e-03 9.31839567e-04]
 [1.04827592e-04 3.51441482e-03 2.05297738e-04 4.02499257e-04]
 [3.87858024e-04 7.99769647e-03 5.85235056e-04 5.91762768e-05]
 [1.26320704e-04 1.84382247e-03 5.84039809e-06 3.39594415e-05]
 [1.38777383e-02 1.73185685e-03 2.00942524e-03 6.09587098e-04]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [2.71560631e-02 1.22320626e-10 1.33757456e-03 1.99071938e-05]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [1.78146584e-03 2.79410650e-03 2.46060202e-03 2.91133741e-02]
 [3.38610920e-03 5.83086773e-03 6.36602543e-02 3.06886333e-03]
 [6.92414686e-03 1.09936351e-01 3.59391950e-03 2.88692781e-04]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [3.52018257e-03 1.28705679e-01 1.44792346e-02 1.61855613e-02]
 [3.28335009e-03 3.20059102e-02 4.79574003e-03 3.53480997e-01]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]]
```

Q-Learning

```
laphisboy@laphisboy-900X3K:~/RL_fall/fall_week_1$ python3 frozen_lake_q.py
Score over time: 0.1205
Final Q-Table Values
[[4.30972397e-02 8.68329874e-02 5.16130150e-03 1.39912677e-02]
 [8.41989861e-04 1.11620981e-03 2.78075826e-04 4.70342591e-02]
 [1.04278780e-03 6.54248387e-03 1.70384432e-03 1.56023476e-03]
 [1.34461696e-03 3.82147058e-05 1.81494164e-04 1.63379968e-03]
 [4.92277658e-02 1.55019952e-03 1.01424568e-02 3.88767743e-03]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [5.30290511e-07 6.58845984e-07 8.85605526e-04 2.62413302e-06]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [1.11079066e-03 1.84759192e-02 3.67341763e-02 3.01863580e-01]
 [4.05955922e-04 5.21394469e-01 3.49797214e-02 5.06608931e-02]
 [1.12056118e-01 2.71397368e-06 8.23377276e-06 2.83359708e-02]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [2.11084497e-02 3.40265637e-01 8.54451439e-01 2.59187179e-03]
 [1.73027131e-01 9.89435350e-01 1.48444394e-01 4.63286997e-01]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]]
```
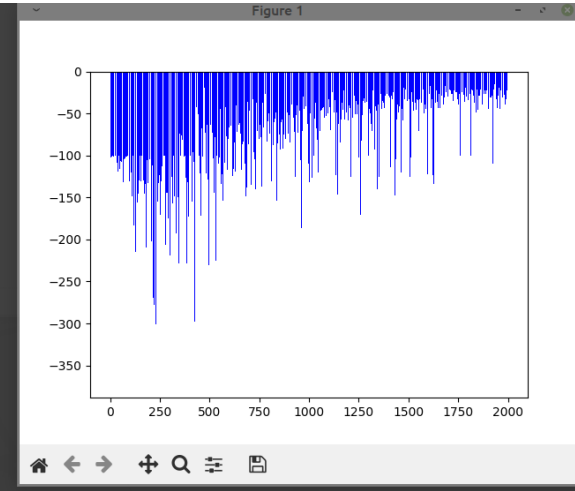
# Why SARSA rather than Q-learning?



SARSA

```
laphisboy@laphisboy-900X3K:~/RL_fall/fall_week_1$ python3 cliff_SARSA.py
Success rate: -78.4395
Final Q-Table Values
[[[ -6.57224026  -6.56897259  -6.57255013  -6.57016348]
  [ -6.40792615  -6.40794207  -6.41262707  -6.40722306]
  [ -6.19563068  -6.19150736  -6.19055137  -6.19008469]
  [ -5.94812063  -5.9433751   -5.94397375  -5.94298732]
  [ -5.66411456  -5.65894466  -5.65843276  -5.66155133]
  [ -5.35078124  -5.34704226  -5.34938056  -5.35235109]
  [ -4.99885228  -4.99516469  -4.99846572  -5.00306954]
  [ -4.600921    -4.59841477  -4.59874203  -4.5972695 ]
  [ -4.152033    -4.15241057  -4.154525    -4.15715972]
  [ -3.66646721  -3.66500498  -3.66553048  -3.66373331]
  [ -3.15142676  -3.14236339  -3.14232914  -3.14794916]
  [ -2.62328346  -2.62218288  -2.6148902   -2.62865687]]

 [[ -6.7159586   -6.71529783  -6.71919877  -6.71900845]
  [ -6.47374734  -6.47359043  -6.84136547  -6.47542736]
  [ -6.21253041  -6.21002613  -6.21058457  -6.21021332]
  [ -5.9561723   -5.95173097  -6.02177763  -5.95312993]
```

Q-Learning

```
laphisboy@laphisboy-900X3K:~/RL_fall/fall_week_1$ python3 cliff_q.py
Success rate: -90.979
Final Q-Table Values
[[[ -5.65282542  -5.65361114  -5.65574407  -5.65471412]
  [ -5.53672858  -5.5375399   -5.54061073  -5.53960431]
  [ -5.35812053  -5.35881434  -5.3613277   -5.36680103]
  [ -5.14629029  -5.14456665  -5.14268356  -5.14482495]
  [ -4.89721249  -4.89824816  -4.8961145   -4.89803612]
  [ -4.62018105  -4.61822491  -4.61625788  -4.62235642]
  [ -4.30511984  -4.30426263  -4.30310824  -4.30619001]
  [ -3.95317341  -3.95352794  -3.95746675  -3.95710484]
  [ -3.56627702  -3.56746228  -3.56728727  -3.57233242]
  [ -3.15481121  -3.15248086  -3.15664835  -3.15322673]
  [ -2.73066757  -2.72614154  -2.72950243  -2.7296559 ]
  [ -2.35850521  -2.35794923  -2.35953948  -2.3586992 ]]

 [[ -5.76232745  -5.76456868  -5.76388702  -5.76406465]
  [ -5.60076494  -5.60036745  -5.60415485  -5.60290239]
  [ -5.39606477  -5.39335139  -5.39800854  -5.39444687]
  [ -5.16188336  -5.15905796  -5.16290388  -5.16118663]
```
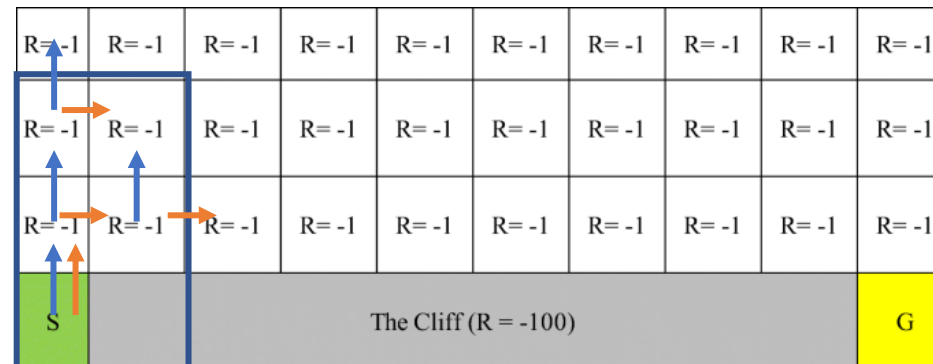
# Why SARSA rather than Q-learning?

```
21    # action = [up, right, down, left]
```

## SARSA

```
55   Q Values near starting position
56   START
57   [[ -5.65773872 -73.19532831  -6.58676967  -6.09652006]]
58   Above START
59   [[ -5.4619545   -5.60901178  -5.85142828  -5.46444349]
60    [ -4.6808752   -4.90508453 -44.73165228  -4.67900451]]
61   Above above START
62   [[-5.18131607 -5.18191879 -5.17902335 -5.18519691]
63    [-4.94702408 -4.94652235 -5.08051739 -4.95036558]]
```

## Q-Learning

```
55   Q Values near starting position
56   START
57   [[ -4.47977623 -87.76072562  -4.48266823  -4.47962649]]
58   Above START
59   [[ -4.12097592  -4.12070405  -4.11925462  -4.11693268]
60    [ -3.8334134   -3.83093316 -57.86657778  -3.83874233]]
61   Above above START
62   [[-3.90748219  3.90242961 -3.90697297 -3.90700077]
63    [-3.75335178 -3.75719326 -3.75394856 -3.75815594]]
```
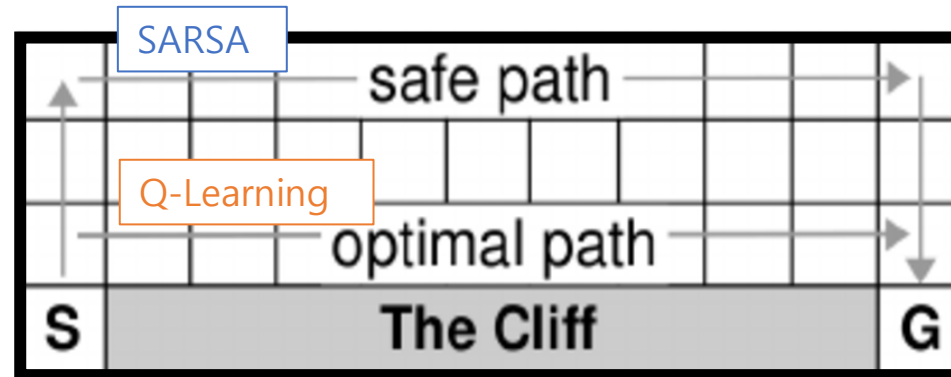
# Why SARSA rather than Q-learning?

```
21    # action = [up, right, down, left]
```



SARSA

safe path

Q-Learning

optimal path

S          The Cliff          G

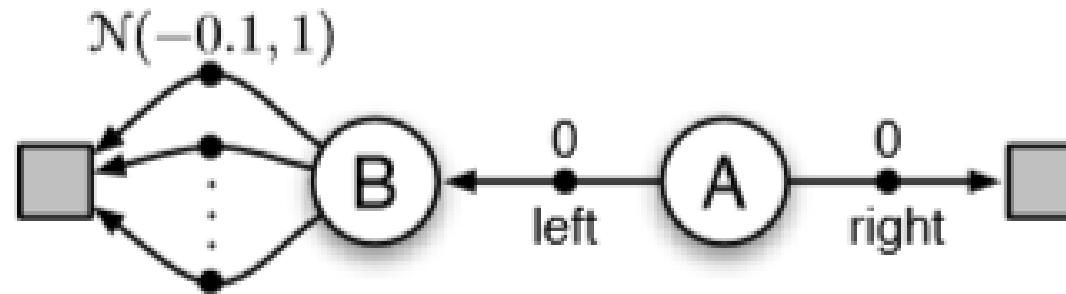| R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 |
|------|------|------|------|------|------|------|------|------|------|
| R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 |
| R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 | R=-1 |
| S | | The Cliff (R = -100) | | | | | | | G |

SARSA → Q-Learning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t))$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

Leads to positive bias

= Maximization Bias

# Maximization Bias



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t))$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

Q. Both have maximization bias right?

# Maximization Bias



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t))$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

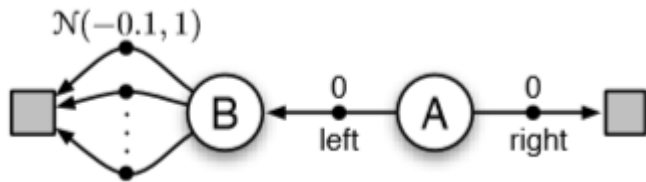Q. Both have maximization bias right?

A. Yes because policy is updated with action with maximum Q-value = optimism in e-greedy policy

B. But SARSA is less optimistic while evaluation of given action

Q-Learning → Double Q-Learning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q_2 \left( S_{t+1}, \arg\max_a Q_1(S_{t+1}, a) \right) - Q_1(S_t, A_t) \right]$$



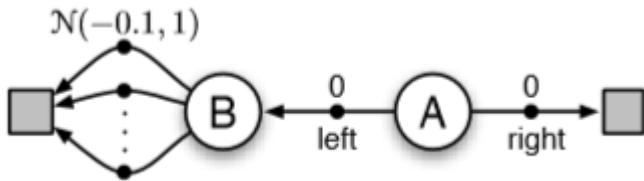Q7. What is double Q-learning? And how does it overcome maximization bias?

# Double Q learning?

---

1: Initialize $Q_1(s, a)$ and $Q_2(s, a), \forall s \in S, a \in A \ t = 0$, initial state $s_t = s_0$

2: **loop**

3:    Select $a_t$ using $\epsilon$-greedy $\pi(s) = \arg\max_a Q_1(s_t, a) + Q_2(s_t, a)$

4:    Observe $(r_t, s_{t+1})$

5:    **if** (with 0.5 probability) **then**

6:       $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha(r_t + \gamma \max_a Q_2(s_{t+1}, a)$

7:    **else**

8:       $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha(r_t + \gamma \max_a Q_1(s_{t+1}, a)$

9:    **end if**

10:    $t = t + 1$

11: **end loop**

---

Q-Learning → Double Q-Learning

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha\left[R_{t+1} + \gamma Q_2\left(S_{t+1}, \arg\max_a Q_1(S_{t+1}, a)\right) - Q_1(S_t, A_t)\right]$$



Q7. What is double Q-learning? And how does it overcome maximization bias?

한 번 왼쪽 action 에서 0보다 큰 값이 나오면 지속적으로 argmax가 left action이 되는 문제가 생김 (true action(?) = right)

Double Q-learning 은 해결이 아닌 완화 method 으로 이해…

# 참고

Base code       https://github.com/hunkim/ReinforcementZeroToAll

See how SARSA is implemented       https://github.com/rlcode
/reinforcement-
learning/tree/master/1-
grid-world

Cliff environment       https://github.com/podondra/gym-gridworlds

Difference between those two       https://tcnguyen.github.io/reinforcement_learning/sarsa_vs_q_learning.html

코드

https://github.com/laphisboy/RL_fall/tree/master/fall_week_1