

RL homework 2

Due date: 26 February 2018, 23:55am

Name: Yifan Shen

ID: 15015762

How to submit

When you have completed the exercises and everything has finished running, click on 'File' in the menu-bar and then 'Download .ipynb'. This file must be submitted to Moodle named as **<student_id>_ucldm_rl2.ipynb** before the deadline above.

Do not forget to include the PDF version on the .zip submitted in Moodle.

Also send a sharable link to the notebook at the following email: **uc1.coursework.submit@gmail.com**

Context

In this assignment, we will take a first look at learning decisions from data. For this, we will use the multi-armed bandit framework.

Background reading

- Sutton and Barto (2018), Chapters 3 - 6

The Assignment

Objectives

You will use Python to implement several reinforcement learning algorithms.

You will then run these algorithms on a few problems, to understand their properties.

Setup

Import Useful Libraries

In [0]:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from collections import namedtuple
```

Set options

In [0]:

```
1 np.set_printoptions(precision=3, suppress=1)
2 plt.style.use('seaborn-notebook')
```

A grid world

In [0]:

```

1 class Grid(object):
2
3     def __init__(self, noisy=False):
4         # -1: wall
5         # 0: empty, episode continues
6         # other: number indicates reward, episode will terminate
7         self._layout = np.array([
8             [-1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1],
9             [-1, 0, 0, 0, 0, 0, 0, 0, 0, 0, -1],
10            [-1, 0, 0, 0, -1, -1, -1, 0, 0, 0, -1],
11            [-1, 0, 0, 0, -1, -1, -1, 0, 10, 0, -1],
12            [-1, 0, 0, 0, -1, -1, -1, 0, 0, 0, -1],
13            [-1, 0, 0, 0, 0, 0, 0, 0, 0, 0, -1],
14            [-1, 0, 0, 0, 0, 0, 0, 0, 0, 0, -1],
15            [-1, 0, 0, 0, 0, 0, 0, 0, 0, 0, -1],
16            [-1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1]
17        ])
18        self._start_state = (2, 2)
19        self._state = self._start_state
20        self._number_of_states = np.prod(np.shape(self._layout))
21        self._noisy = noisy
22
23    @property
24    def number_of_states(self):
25        return self._number_of_states
26
27    def plot_grid(self):
28        plt.figure(figsize=(4, 4))
29        plt.imshow(self._layout > -1, interpolation="nearest", cmap='pink')
30        ax = plt.gca()
31        ax.grid(0)
32        plt.xticks([])
33        plt.yticks([])
34        plt.title("The grid")
35        plt.text(2, 2, r"$\mathbf{S}$", ha='center', va='center')
36        plt.text(8, 3, r"$\mathbf{G}$", ha='center', va='center')
37        h, w = self._layout.shape
38        for y in range(h-1):
39            plt.plot([-0.5, w-0.5], [y+0.5, y+0.5], '-k', lw=2)
40        for x in range(w-1):
41            plt.plot([x+0.5, x+0.5], [-0.5, h-0.5], '-k', lw=2)
42
43
44    def get_obs(self):
45        y, x = self._state
46        return y*self._layout.shape[1] + x
47
48    def obs_to_state(obs):
49        x = obs % self._layout.shape[1]
50        y = obs // self._layout.shape[1]
51        s = np.copy(self._layout)
52        s[y, x] = 4
53        return s
54
55    def step(self, action):
56        y, x = self._state
57
58        if action == 0: # up
59            new_state = (y - 1, x)

```

```
60     elif action == 1: # right
61         new_state = (y, x + 1)
62     elif action == 2: # down
63         new_state = (y + 1, x)
64     elif action == 3: # left
65         new_state = (y, x - 1)
66     else:
67         raise ValueError("Invalid action: {} is not 0, 1, 2, or 3.".format(action))
68
69     new_y, new_x = new_state
70     if self._layout[new_y, new_x] == -1: # wall
71         reward = -5.
72         discount = 0.9
73         new_state = (y, x)
74     elif self._layout[new_y, new_x] == 0: # empty cell
75         reward = 0.
76         discount = 0.9
77     else: # a goal
78         reward = self._layout[new_y, new_x]
79         discount = 0.
80         new_state = self._start_state
81     if self._noisy:
82         width = self._layout.shape[1]
83         reward += 2*np.random.normal(0, width - new_x + new_y)
84
85     self._state = new_state
86
87     return reward, discount, self.get_obs()
```

Helper functions

In [0]:

```

1 def run_experiment(env, agent, number_of_steps):
2     mean_reward = 0.
3     try:
4         action = agent.initial_action()
5     except AttributeError:
6         action = 0
7     for i in range(number_of_steps):
8         reward, discount, next_state = grid.step(action)
9         action = agent.step(reward, discount, next_state)
10        mean_reward += (reward - mean_reward)/(i + 1.)
11
12    return mean_reward
13
14 map_from_action_to_subplot = lambda a: (2, 6, 8, 4)[a]
15 map_from_action_to_name = lambda a: ("up", "right", "down", "left")[a]
16
17 def plot_values(values, colormap='pink', vmin=0, vmax=10):
18     plt.imshow(values, interpolation="nearest", cmap=colormap, vmin=vmin, vmax=vmax)
19     plt.yticks([])
20     plt.xticks([])
21     plt.colorbar(ticks=[vmin, vmax])
22
23 def plot_action_values(action_values, vmin=0, vmax=10):
24     q = action_values
25     fig = plt.figure(figsize=(8, 8))
26     fig.subplots_adjust(wspace=0.3, hspace=0.3)
27     for a in [0, 1, 2, 3]:
28         plt.subplot(3, 3, map_from_action_to_subplot(a))
29         plot_values(q[... , a], vmin=vmin, vmax=vmax)
30         action_name = map_from_action_to_name(a)
31         plt.title(r"$q(s, \mathrm{" + action_name + r"})$")
32
33     plt.subplot(3, 3, 5)
34     v = 0.9 * np.max(q, axis=-1) + 0.1 * np.mean(q, axis=-1)
35     plot_values(v, colormap='summer', vmin=vmin, vmax=vmax)
36     plt.title("$v(s)$")
37
38
39 def plot_rewards(xs, rewards, color):
40     mean = np.mean(rewards, axis=0)
41     p90 = np.percentile(rewards, 90, axis=0)
42     p10 = np.percentile(rewards, 10, axis=0)
43     plt.plot(xs, mean, color=color, alpha=0.6)
44     plt.fill_between(xs, p90, p10, color=color, alpha=0.3)
45
46
47 def parameter_study(parameter_values, parameter_name,
48                     agent_constructor, env_constructor, color, repetitions=10, number_of_steps=int
49                     mean_rewards = np.zeros((repetitions, len(parameter_values)))
50                     greedy_rewards = np.zeros((repetitions, len(parameter_values)))
51     for rep in range(repetitions):
52         for i, p in enumerate(parameter_values):
53             env = env_constructor()
54             agent = agent_constructor()
55             if 'eps' in parameter_name:
56                 agent.set_epsilon(p)
57             elif 'alpha' in parameter_name:
58                 agent._step_size = p
59             else:

```

```

60         raise NameError("Unknown parameter_name: {}".format(parameter_name))
61     mean_rewards[rep, i] = run_experiment(grid, agent, number_of_steps)
62     agent.set_epsilon(0.)
63     agent._step_size = 0.
64     greedy_rewards[rep, i] = run_experiment(grid, agent, number_of_steps//10)
65     del env
66     del agent
67
68     plt.subplot(1, 2, 1)
69     plot_rewards(parameter_values, mean_rewards, color)
70     plt.yticks([0, 1], [0, 1])
71     # plt.ylim((0, 1.5))
72     plt.ylabel("Average reward over first {} steps".format(number_of_steps), size=12)
73     plt.xlabel(parameter_name, size=12)
74
75     plt.subplot(1, 2, 2)
76     plot_rewards(parameter_values, greedy_rewards, color)
77     plt.yticks([0, 1], [0, 1])
78     # plt.ylim((0, 1.5))
79     plt.ylabel("Final rewards, with greedy policy".format(number_of_steps), size=12)
80     plt.xlabel(parameter_name, size=12)
81
82 def epsilon_greedy(q_values, epsilon):
83     if epsilon < np.random.random():
84         return np.argmax(q_values)
85     else:
86         return np.random.randint(np.array(q_values).shape[-1])

```

Part 1: Implement agents

Each agent, should implement a step function:

step(self, reward, discount, next_observation, ...):

where ... indicates there could be other inputs (discussed below). The step should update the internal values, and return a new action to take.

When the discount is zero ($\text{discount} = \gamma = 0$), then the `next_observation` will be the initial observation of the next episode. One shouldn't bootstrap on the value of this state, which can simply be guaranteed when using " $\gamma \cdot v(\text{next_observation})$ " (for whatever definition of v is appropriate) in the update, because $\gamma = 0$. So, the end of an episode can be seamlessly handled with the same step function.

__init__(self, number_of_actions, number_of_states, initial_observation):

The constructor will provide the agent the number of actions, number of states, and the initial observation. You can get the initial observation by first instantiating an environment, using `grid = Grid()`, and then calling `grid.get_obs()`.

In this assignment, observations will be states in the environment, so the agent state, environment state, and observation will overlap, and we will use the word `state` interchangeably with `observation`.

All agents should be in pure Python - so you cannot use TensorFlow to, e.g., compute gradients. Using numpy is fine.

A note on the initial action

Normally, you would also have to implement a method that gives the initial action, based on the initial state. In our experiments the helper functions above will just use the action 0 (which corresponds to up) as initial action, so that otherwise we do not have to worry about this. Note that this initial action is only executed once, and the beginning of the first episode---not at the beginning of each episode.

Some algorithms (Q-learning, Sarsa) need to remember the last action in order to update its value when they see the next state. In the `__init__`, make sure you set the initial action to zero, e.g.,

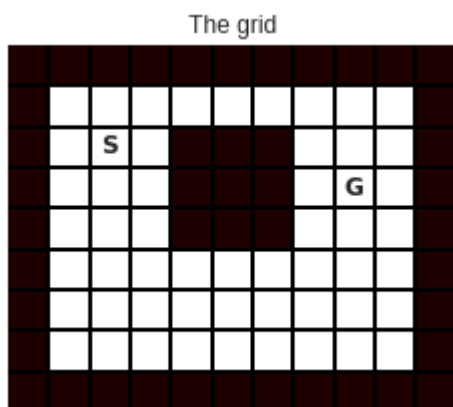
```
def __init__(...):
    (...)
    self._last_action = 0
    (...)
```

The grid

The cell below shows the Grid environment that we will use. Here S indicates the start state and G indicates the goal. The agent has four possible actions: up, right, down, and left. Rewards are: -5 for bumping into a wall, +10 for reaching the goal, and 0 otherwise. The episode ends when the agent reaches the goal, and otherwise continues. The discount, on continuing steps, is $\gamma = 0.9$. Feel free to reference the implementation of the Grid above, under the header "a grid world".

In [0]:

```
1 grid = Grid()
2 grid.plot_grid()
```



Random agent

In [0]:

```
1 # For reference: here is a random agent
2 class Random(object):
3
4     def __init__(self, number_of_actions, number_of_states, initial_state):
5         self._number_of_actions = number_of_actions
6
7     def step(self, reward, discount, next_state):
8         next_action = np.random.randint(number_of_actions)
9         return next_action
```

Agent 1: TD learning

[5 pts] Implement an agent that behaves randomly, but that *on-policy* estimates state values $v(s)$, using one-step TD learning with a step size $\alpha = 0.1$.

Also implement `get_values(self)` that returns the vector of all state values (one value per state).

You should be able to use the `__init__` as provided below, so you just have to implement `get_values` and `step`. We store the initial state in the constructor because you need its value on the first step in order to compute the TD error when the first transition has occurred. Hint: in the `step` you similarly will want to store the previous state to be able to compute the next TD error on the next step.

In [0]:

```

1 class RandomTD(object):
2
3     def __init__(self, number_of_states, number_of_actions, initial_state, step_size):
4         self._values = np.zeros(number_of_states)
5         self._state = initial_state
6         self._number_of_actions = number_of_actions
7         self._step_size = step_size
8
9     def get_values(self):
10         return self._values
11
12     def step(self, r, g, s):
13         # r: reward
14         # g: discount rate
15         # s: state
16         self._values[self._state] += \
17             self._step_size * (r + g * self._values[s] - self._values[self._state])
18         self._state = s
19         return np.random.randint(self._number_of_actions)

```

Run the next cell to run the RandomTD agent on a grid world.

In [39]:

```

1 # DO NOT MODIFY THIS CELL
2 agent = RandomTD(grid._layout.size, 4, grid.get_obs())
3 run_experiment(grid, agent, int(1e5))
4 v = agent.get_values()
5 plot_values(v.reshape(grid._layout.shape), colormap="hot", vmin=-10, vmax=5)

```



If everything worked as expected, the plot above will show the estimated state values under the random policy. This includes values for unreachable states --- on the walls and on the goal (we never actually reach the goal -- - rather, the episode terminates on the transition to the goal. The values on the walls and goal are, and will always remain, zero (shown in orange above).

Policy iteration

We used TD to do policy evaluation for the random policy on this problem. Consider doing policy improvement, by taking the greedy policy with respect to a one-step look-ahead. For this, you may assume we have a true model, so for each state the policy would for each action look at the value of the resulting state, and would then pick the action with the highest state value. You do **not** have to implement this, just answer the following questions.

[5 pts] Would the greedy after one such iteration of policy evaluation and policy improvement be optimal on this problem? Explain (in one or two sentences) why or why not.

Answer: No, one iteration will not reach optimal. One iteration will lead to a better policy but not the optimal.

[5 pts] If we repeat the process over and over again, and repeatedly evaluate the greedy policy and then perform another improvement step, would then the policy eventually become optimal? Explain (in one or two sentences) why or why not.

Answer: Yes, multiple iterations will reach optimal. From policy improvement theorem (Sutton's book), if we perform another improvement step, we will obtain a better policy. Hence if we repeat the process (policy iteration), we will finally reach the optimal policy.

Agent 2: Q-learning

[10 pts] Implement an agent that uses **Q-learning** to learn action values. In addition, the agent should act according to an ϵ -greedy policy with respect to its action values.

[10 pts] Include an option to use **Double Q-learning**, with a `double` boolean flag in the `init`. If `double=False` the agent should perform Q-learning. If `double=True` the agent should perform Double Q-learning. Note that we then need two action-value functions.

[10 pts] Include an option to use **Sarsa** instead of Q-learning, in the `step`.

[15 pts] Generalize to **General Q-learning**, where the `__init__` takes a functions `target_policy` and `behaviour_policy`. The function `behaviour_policy(action_values)` should map `action_values` to a single action. For instance, the random policy can be implemented as:

```
def behaviour_policy(action_values):
    return np.random.randint(len(action_values))
```

We will typically just use ϵ -greedy, for instance:

```
def behaviour_policy(action_values):
    return epsilon_greedy(action_values, epsilon=0.1)
```

The target policy is defined by a function `target_policy(action_values, action)`, which should return **a vector** with one probability per action. The `action` argument is used to be able to do Sarsa: in addition to the action values, the function will also get the action as selected by the behaviour so that it can return a one hot vector for just the selected action in the Sarsa case. So, the target policy for Sarsa would look like this:

```
def one_hot(index, max_index):
    np.eye(max_index)[index]

def target_policy(action_values, action):
    return one_hot(action)
```

As another example, a random target policy is:

```
def target_policy(action_values, unused_action):
    number_of_actions = len(action_values)
    return np.ones((number_of_actions,))/number_of_actions
```

Note that **double learning** can be combined with General Q-learning, but is separate. So, the `double` flag in the `init` remains. E.g., when the target policy is the Sarsa policy described above and `double=True`, the algorithm should implement **double Sarsa**.

Note also that if (or when) you have implemented General Q-learning, this algorithm encompasses all previous algorithms, so you only need this one algorithm with, as its interface the two functions

```
__init__(self, number_of_states, number_of_actions, initial_state, target_policy,
behaviour_policy, double, step_size=0.1)
```

and

```
step(self, reward, discount, next_state)
```

We will mostly use `step_size=0.1`, so make that the default, but allow it to change when it is fed in as an argument.

If you do not success in implementing General Q-learning, try to implement at least Sarsa, Q-learning, and Double Q-learning, to be able to answer questions below.

In [0]:

```

1 class GeneralQ(object):
2
3     def __init__(self, number_of_states, number_of_actions, initial_state, \
4                 target_policy, behaviour_policy, double, step_size=0.1):
5         self._q = np.zeros((number_of_states, number_of_actions))
6         if double:
7             self._q2 = np.zeros((number_of_states, number_of_actions))
8         self._s = initial_state
9         self._number_of_actions = number_of_actions
10        self._step_size = step_size
11        self._behaviour_policy = behaviour_policy
12        self._target_policy = target_policy
13        self._double = double
14        self._last_action = 0
15
16    @property
17    def q_values(self):
18        if self._double:
19            return (self._q + self._q2)/2
20        else:
21            return self._q
22
23    def step(self, r, g, s):
24
25        if self._double:
26            # execute double Q-learning
27            # choose a from s
28            a = self._behaviour_policy(self.q_values[s,:]) \
29            # integer 0 to 3, chosen by epsilon_greedy
30
31            # s: S dash, new state
32            # self._s: Capital S, current state
33            if np.random.rand() < .5: # with 0.5 probability
34                action_one_hot = self._target_policy(self._q[s,:), a) \
35                # one-hot vector, chosen max
36                self._q[self._s, self._last_action] += \
37                self._step_size * (r + g * np.dot(self._q2[s,:), action_one_hot) \
38                - self._q[self._s, self._last_action] )
39            else: # with 0.5 probability
40                action_one_hot = self._target_policy(self._q2[s,:), a) \
41                # one-hot vector, chosen max
42                self._q2[self._s, self._last_action] += \
43                self._step_size * (r + g * np.dot(self._q[s,:), action_one_hot) \
44                - self._q2[self._s, self._last_action] )
45        else:
46            # execute Q-learning
47            a = self._behaviour_policy(self.q_values[s,:])
48            action_one_hot = self._target_policy(self._q[s,:), a) \
49            # one-hot vector, chosen max
50            self._q[self._s, self._last_action] += \
51            self._step_size * (r + g * np.dot(self._q[s,:), action_one_hot) \
52            - self._q[self._s, self._last_action] )
53        # update a and s
54        self._last_action = a
55        self._s = s
56        return self._last_action

```

Part 2: Analyse Results

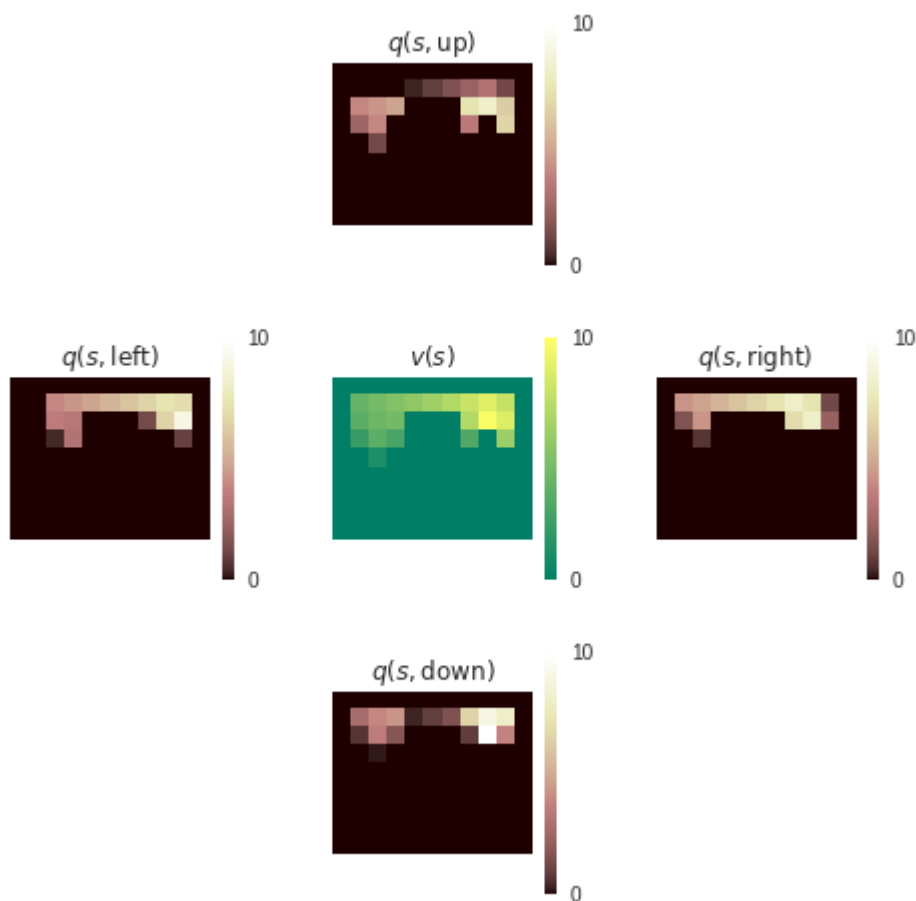
Run the cells below to train a Q-learning and a SARSA agent and generate plots.

This trains the agents the Grid problem with a step size of 0.1 and an epsilon of 0.1.

The plots below will show action values for each of the actions, as well as a state value defined by $v(s) = \sum_a \pi(as)q(s, a)$.

In [41]:

```
1 # Q-learning
2 grid = Grid()
3 def target_policy(q, a):
4     return np.eye(len(q))[np.argmax(q)]
5 def behaviour_policy(q):
6     return epsilon_greedy(q, 0.1)
7 agent = GeneralQ(grid._layout.size, 4, grid.get_obs(),
8                 target_policy, behaviour_policy, double=False)
9 run_experiment(grid, agent, int(1e5))
10 q = agent.q_values.reshape(grid._layout.shape + (4,))
11 plot_action_values(q)
```

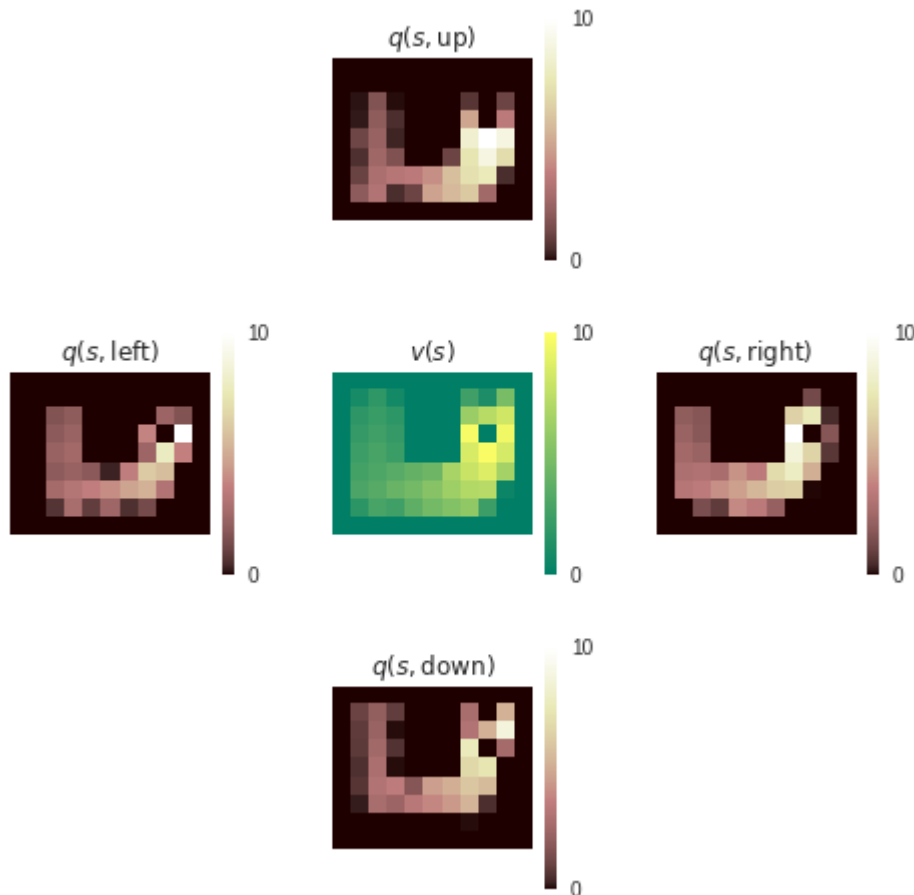


In [42]:

```

1 # Sarsa
2 grid = Grid()
3 def target_policy(q, a):
4     return np.eye(len(q))[a]
5 def behaviour_policy(q):
6     return epsilon_greedy(q, 0.1)
7 agent = GeneralQ(grid._layout.size, 4, grid.get_obs(),
8                 target_policy, behaviour_policy, double=False)
9 run_experiment(grid, agent, int(1e5))
10 q = agent.q_values.reshape(grid._layout.shape + (4,))
11 plot_action_values(q)

```



Questions

Consider the greedy policy with respect to the estimated values

[10 pts] How do the policies found by Q-learning and Sarsa differ? (Explain qualitatively how the behaviour differs in one or two sentences.)

Answer: Policies found by Q-learning tends to move above the square (has less exploration, and less steps), but while policies found by Sarsa tends to move below the square (has more exploration and more steps).

[10 pts] Why do the policies differ in this way?

Answer: In the Q-learning control policy, the action to take is chosen by having the highest action value. It follows the optimal control strategy, therefore the action values will converge such that the best path is above the square. However, in the Sarsa control policy, the agent's actual control strategy is taken into account when learning. Therefore, the agent has learned from time to time that avoiding going

above the square to avoid losses (it is more likely to hit wall and get losses when travelling above the square).

[10 pts] Which greedy policy is better, in terms of actual value?

Answer: Sarsa performs better as it explores more states and hence avoid hitting the wall. Hence having the higher actual value.

Noisy environments

We will now compare Q-learning and Double Q-learning on a noisy version of the environment.

In the noisy version, a zero-mean Gaussian is added to all rewards. The variance of this noise is higher the further to the left you go, and the further down (so away from the goal).

Run the cell below to run 20 repetitions of the experiment that runs Q-learning and Double Q-learning on this noisy domain.

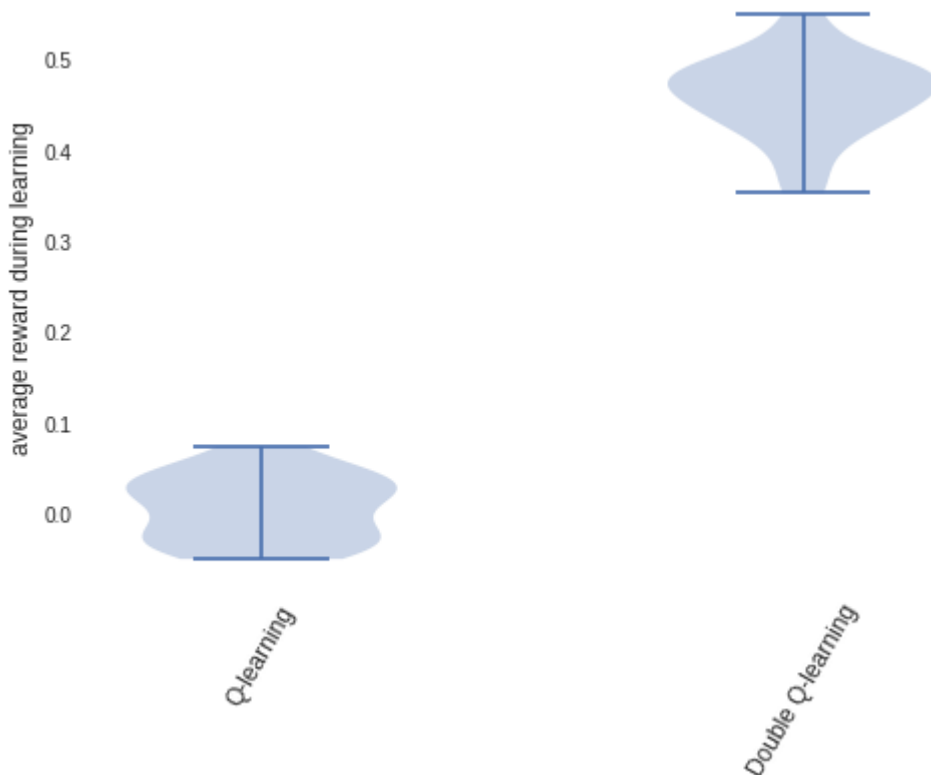
In [43]:

```

1 def target_policy(q, a):
2     max_q = np.max(q)
3     pi = np.array([1. if qi == max_q else 0. for qi in q])
4     return pi / sum(pi)
5 def behaviour_policy(q):
6     return epsilon_greedy(q, 0.1)
7 mean_reward_q_learning = []
8 mean_reward_double_q_learning = []
9 for _ in range(20):
10     grid = Grid(noisy=True)
11     q_agent = GeneralQ(grid._layout.size, 4, grid.get_obs(),
12                        target_policy, behaviour_policy, double=False, step_size=0)
13     dq_agent = GeneralQ(grid._layout.size, 4, grid.get_obs(),
14                        target_policy, behaviour_policy, double=True, step_size=0)
15     mean_reward_q_learning.append(run_experiment(grid, q_agent, int(2e5)))
16     mean_reward_double_q_learning.append(run_experiment(grid, dq_agent, int(2e5)))
17 plt.violinplot([mean_reward_q_learning, mean_reward_double_q_learning])
18 plt.xticks([1, 2], ["Q-learning", "Double Q-learning"], rotation=60, size=12)
19 plt.ylabel("average reward during learning", size=12)
20 ax = plt.gca()
21 ax.set_axis_bgcolor('white')
22 ax.grid(0)

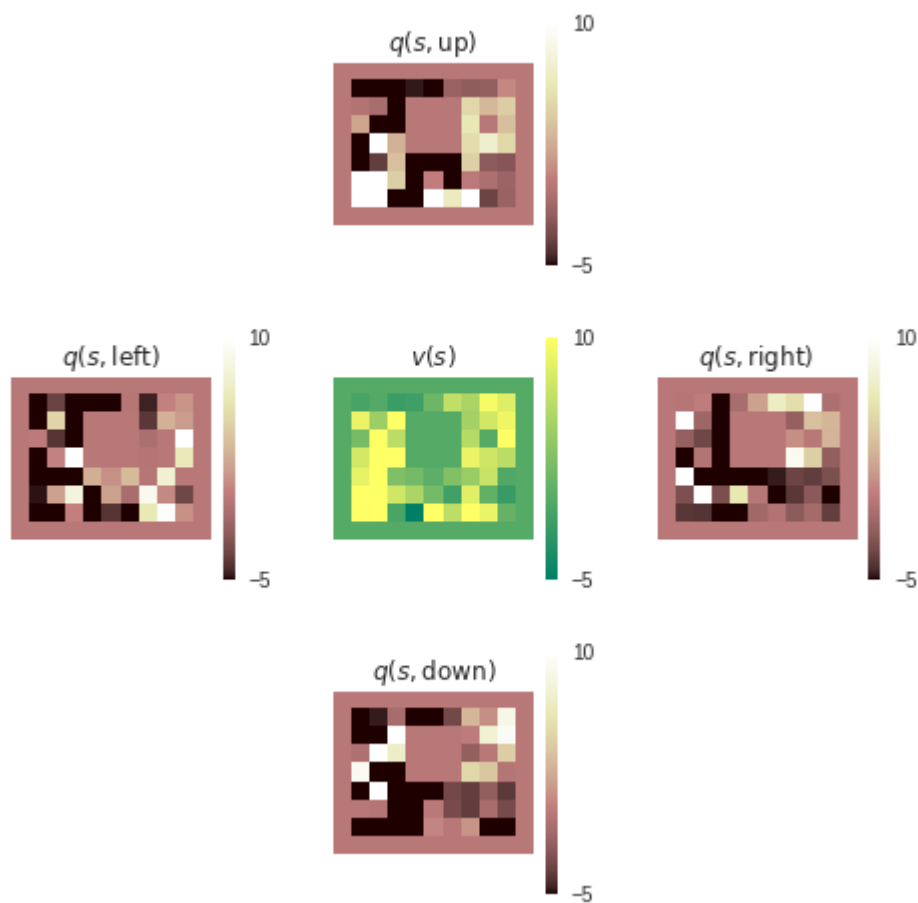
```

/usr/local/lib/python2.7/dist-packages/ipykernel_launcher.py:21: MatplotlibDeprecationWarning: The set_axis_bgcolor function was deprecated in version 2.0. Use set_facecolor instead.



In [44]:

```
1 # single Q-learning
2 q = q_agent.q_values.reshape(grid._layout.shape + (4,))
3 plot_action_values(q, vmin=-5)
```

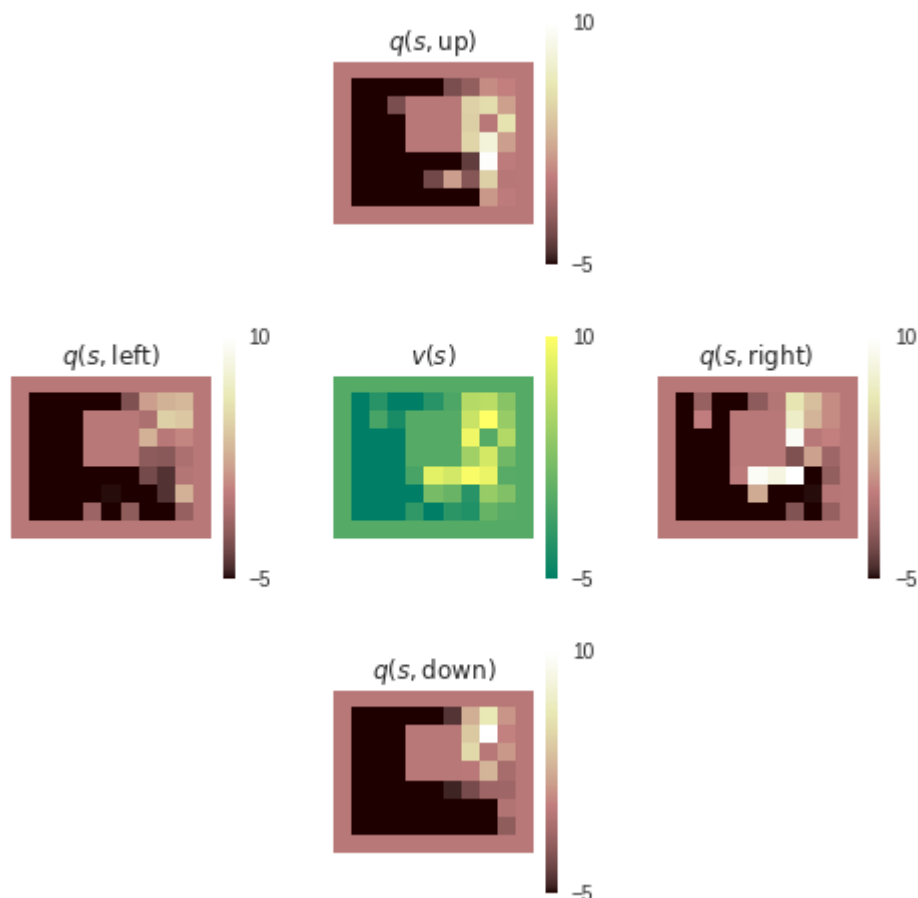


In [45]:

```

1 # Double Q-learning
2 q = dq_agent.q_values.reshape(grid._layout.shape + (4,))
3 plot_action_values(q, vmin=-5)

```



The plots above show 1) the distributions of average rewards (over all learning steps) over the 20 experiments per algorithm, 2) the action values for Q-learning, and 3) the action values for Double Q-learning.

[10 pts] Explain why Double Q-learning has a higher average reward. Use at most four sentences, and discuss at least a) the dynamics of the algorithm, b) how this affects behaviour, and c) why the resulting behaviour yields higher rewards for Double Q-learning than for Q-learning.

(Feel free to experiment to gain more intuition, if this is helpful. Especially the action value plots can be quite noisy, and therefore hard to interpret.)

Answer: In double Q-learning, we have two Q maxtrices which stores value, actions are taken and update each other with 0.5 probability. When using "single" Q-learning under a noisy environment, we might have a higher positive bias (overesimate values), while using Double Q-learning, we averages out the noise and hence will have less bias. Therefore, double Q-learning performs better under a noisy environment, as it averages the noisy return, hence obtain a better estimation of the true value. leading to a higher reward.

In [0]:

1

