

Introduction to Reinforcement Learning

OCNC 2019

```
In [1]: 1 import numpy as np
        2 import matplotlib.pyplot as plt
        3 %matplotlib inline
```

Classes for minimum environment and agent

```
In [2]: 1 class Environment:
        2     """Class for a reinforcement learning environment"""
        3
        4     def __init__(self, nstate=3, naction=2):
        5         """Create a new environment"""
        6         self.Ns = nstate # number of states
        7         self.Na = naction # number of actions
        8
        9     def start(self):
       10         """start an episode"""
       11         # randomly pick a state
       12         self.state = np.random.randint(self.Ns)
       13         return self.state
       14
       15     def step(self, action):
       16         """step forward given an action"""
       17         # random reward
       18         self.reward = np.random.random() # between 0 and 1
       19         # random state transition
       20         self.state = np.random.randint(self.Ns)
       21         return self.reward, self.state
```

In [3]:

```
1 class Agent:
2     """Class for a reinforcement learning agent"""
3
4     def __init__(self, nstate, naction):
5         """Create a new agent"""
6         self.Ns = nstate # number of states
7         self.Na = naction # number of actions
8
9     def start(self, state):
10        """first action, without reward feedback"""
11        # randomly pick an action
12        self.action = np.random.randint(self.Na)
13        return self.action
14
15    def step(self, reward, state):
16        """learn by reward and take an action"""
17        # do nothing for reward
18        # randomly pick an action
19        self.action = np.random.randint(self.Na)
20        return self.action
```

In [4]:

```
1 class RL:
2     """Reinforcement learning by interaction of Environment and Agent"""
3
4     def __init__(self, environment, agent, nstate, naction):
5         """Create the environment and the agent"""
6         self.env = environment(nstate, naction)
7         self.agent = agent(nstate, naction)
8
9     def episode(self, tmax=50):
10        """One episode"""
11        # First contact
12        state = self.env.start()
13        action = self.agent.start(state)
14        # Table of t, r, s, a
15        Trsa = np.zeros((tmax+1,4))
16        Trsa[0,:] = [0, 0, state, action]
17        # Repeat interaction
18        for t in range(1, tmax+1):
19            reward, state = self.env.step(action)
20            action = self.agent.step(reward, state)
21            Trsa[t,:] = [t, reward, state, action]
22        return Trsa
23
24    def run(self, nrun=10, tmax=50):
25        """Multiple runs of episodes"""
26        Return = np.zeros(nrun)
27        for n in range(nrun):
28            r = self.episode(tmax)[:,-1] # reward sequence
29            Return[n] = sum(r)
30        return Return
```

Q learning of Pain-Gain task

In [5]:

```
1 class PainGain(Environment):
2     """Pain-Gain environment """
3
4     def __init__(self, nstate=4, naction=2, gain=6):
5         super().__init__(nstate, naction)
6         # setup the reward function as an array
7         self.R = np.ones((self.Ns, self.Na))
8         self.R[:,1] = -1 # small pains for action 1
9         self.R[0,0] = -gain # large pain
10        self.R[-1,1] = gain # large gain
11
12    def step(self, action):
13        """step by an action"""
14        self.reward = self.R[self.state, action] # reward
15        self.state = (self.state + 2*action-1)%self.Ns # move left
16        return(self.reward, self.state)
```

In [6]:

```

1 class QL(Agent):
2     """Class for a Q-learning agent"""
3
4     def __init__(self, nstate, naction):
5         super().__init__(nstate, naction)
6         # allocate Q table
7         self.Q = np.zeros((nstate, naction))
8         # default parameters
9         self.alpha = 0.1 # learning rate
10        self.beta = 2.0 # inverse temperature
11        self.gamma = 0.9 # discount factor
12
13        def boltzmann(self, q):
14            """Boltzmann selection"""
15            pa = np.exp( self.beta*q) # unnormalized probability
16            pa = pa/sum(pa) # normalize
17            return np.random.choice(self.Na, p=pa)
18
19        def start(self, state):
20            """first action, without reward feedback"""
21            # Boltzmann action selection
22            self.action = self.boltzmann( self.Q[state,:])
23            # remember the state
24            self.state = state
25            return self.action
26
27        def step(self, reward, state):
28            """learn by reward and take an action"""
29            # TD error: self.state/action retains the previous ones
30            delta = reward + self.gamma*max(self.Q[state,:]) - self.Q[s
31            # Update the value for previous state and action
32            self.Q[self.state,self.action] += self.alpha*delta
33            # Boltzmann action selection
34            self.action = self.boltzmann( self.Q[state,:])
35            # remember the state
36            self.state = state
37            return self.action

```

Setup and Run

In [7]:

```

1 # Setup Pain-Gain environment and Q-learning agent
2 pgq = RL(PainGain, QL, 4, 2)
3 pgq.env.R

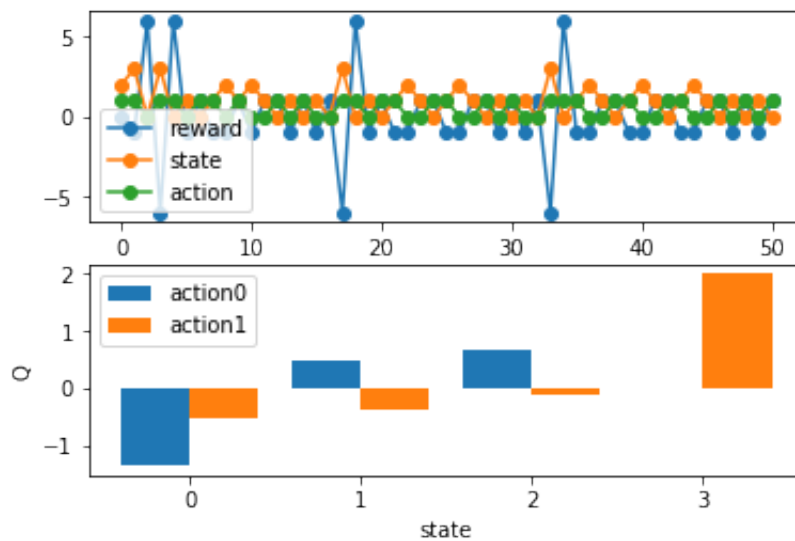
```

```

Out[7]: array([[ -6.,  -1.],
               [  1.,  -1.],
               [  1.,  -1.],
               [  1.,   6.]])

```

```
In [8]: 1 # Run an episode of 50 trials
2 trsa = pgq.episode(50)
3 # plot the time course of state, action, reward,
4 plt.subplot(2,1,1)
5 plt.plot(trsa[:,1:], "o-")
6 plt.legend(["reward", "state", "action"])
7 plt.xlabel("time")
8 # plot the action values
9 plt.subplot(2,1,2)
10 plt.bar(np.arange(pgq.agent.Ns)-0.2, pgq.agent.Q[:,0], 0.4) # act.
11 plt.bar(np.arange(pgq.agent.Ns)+0.2, pgq.agent.Q[:,1], 0.4) # act.
12 plt.legend(["action0", "action1"]);
13 plt.xticks(range(pgq.agent.Ns)); plt.xlabel("state"); plt.ylabel("Q")
```



```
In [9]: 1 # reduce the discount factor
2 pgq.agent.gamma = 0.3
```

```
In [10]: 1 # restore the discount factor
2 pgq.agent.gamma = 0.9
```

```
In [11]: 1 # increase the inverse temperature
2 pgq.agent.beta = 10
```

```
In [12]: 1 # reduce the inverse temperature
2 pgq.agent.beta = 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

