

Day 2. Tabular Q-learning

SAMSUNG AI

Reinforcement Learning

June 17, 2022

Jaeuk Shin, Mingyu Park



CORE
Control + Optimization Research Lab

Tabular Q-learning

With Q, we don't need a policy!

$$\begin{aligned} Q^*(s, a) &= \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \underbrace{\sum_{s' \in S} p(s'|s, a) v^*(s')}_{\text{optimal value of next state}} \\ &= r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) \max_{a' \in A} Q^*(s', a') \end{aligned}$$

- Define the Bellman operator \mathcal{T} for Q-functions by

$$(\mathcal{T}Q)(s, a) := r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) \max_{a' \in A} Q(s', a').$$

Then, it is a monotone contraction mapping.

- Bellman equation:

$$Q = \mathcal{T}Q.$$

Q-Learning (tabular):

Initialize Q ;

Replace model with data!

- Take some action and observe (s, a, s', r) ;
- Set $Q(s, a) \leftarrow (1 - \alpha) \underbrace{Q(s, a)}_{\text{old estimate}} + \alpha \underbrace{[r + \gamma \max_{a'} Q(s', a')]}_{\text{new estimate}};$
- Repeat until convergence;

Can you see the difference?



Tabular Q-learning

MDP definition

```
1  class MyEnv:
2      num_actions = 4
3
4      def __init__(self):
5          pass
6
7      def reset(self):
8          pass
9
10     def step(self, action):
11         pass
```

s_t is kept internally, and is updated in **step** method

sample an initial state $s_0 \sim \rho_0(s)$

agent-env interaction:

$$s_{t+1} \sim p(\cdot | s_t, a_t), \quad r_t = r(s_t, a_t)$$

Tabular Q-learning

Example - Pendulum

```
1  class PendulumEnv(gym.Env):  
2      metadata = {  
3          'render.modes': ['human', 'rgb_array'],  
4          'video.frames_per_second': 30  
5      }  
6  
7      def __init__(self, g=10.0):  
8          self.max_speed = 8  
9          self.max_torque = 2.  
10         self.dt = .05  
11         self.g = g  
12         self.m = 1.  
13         self.l = 1.  
14         self.viewer = None  
15  
16         high = np.array([1., 1., self.max_speed], dtype=np.float32)  
17         self.action_space = spaces.Box(  
18             low=-self.max_torque,  
19             high=self.max_torque, shape=(1,),  
20             dtype=np.float32  
21         )
```

MDP as a Python class(gym.Env)

https://github.com/openai/gym/blob/master/gym/envs/classic_control/pendulum.py

Tabular Q-learning

Example - Pendulum

action input a_t (torque applied at t)

```
1 def step(self, u):  
2     th, thdot = self.state # th := theta
```

state s_t : $s_t = (\theta_t, \dot{\theta}_t)$

```
3  
4     g = self.g  
5     m = self.m  
6     l = self.l  
7     dt = self.dt
```

```
8  
9     u = np.clip(u, -self.max_torque, self.max_torque)[0]  
10    self.last_u = u # for rendering
```

```
11    costs = angle_normalize(th) ** 2 + .1 * thdot ** 2 + .001 * (u ** 2)
```

compute reward $r_t = r(s_t, a_t)$

```
12  
13    newthdot = thdot + (-3 * g / (2 * l) * np.sin(th + np.pi) + 3. / (m * l ** 2) * u) * dt  
14    newth = th + newthdot * dt  
15    newthdot = np.clip(newthdot, -self.max_speed, self.max_speed)
```

compute the next state s_{t+1}

```
16  
17    self.state = np.array([newth, newthdot])  
18    return self._get_obs(), -costs, False, {}
```

return next state s_{t+1} & reward r_t

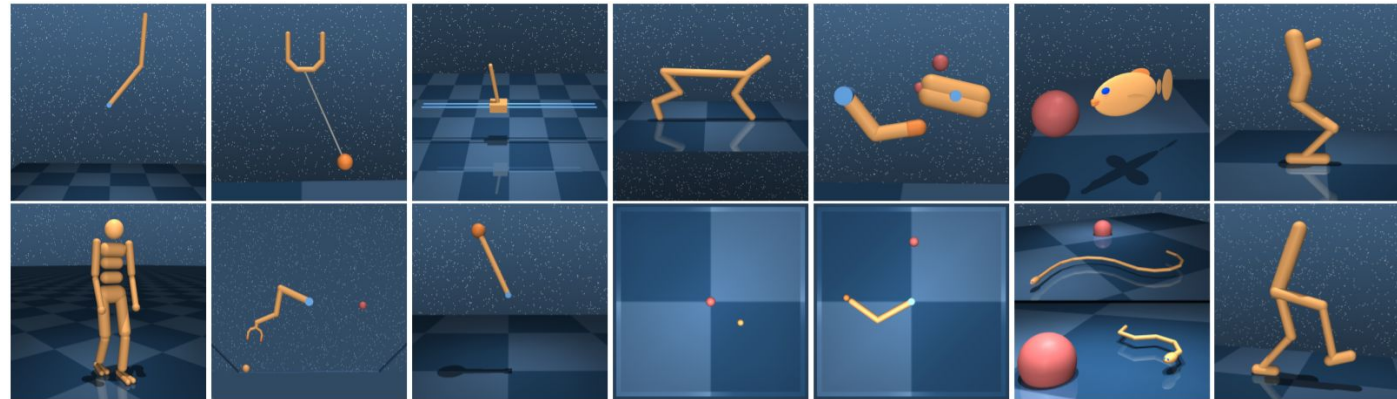
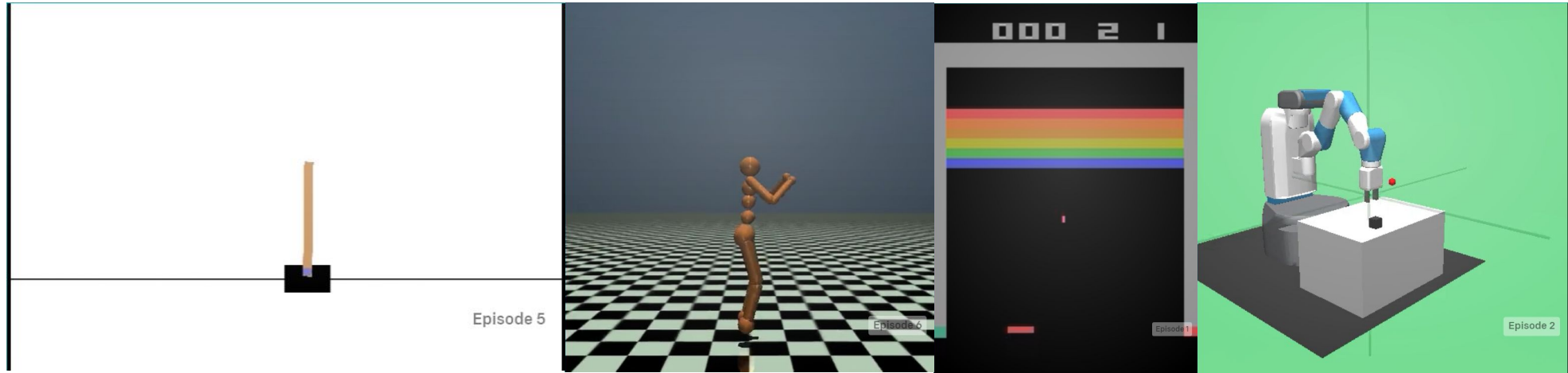
https://github.com/openai/gym/blob/master/gym/envs/classic_control/pendulum.py



CORE
Control + Optimization Research Lab

Tabular Q-learning

More examples - OpenAI Gym, Deepmind Control Suite, etc.



CORE
Control + Optimization Research Lab

Q-learning - Implementation

Algorithm Implementation

```
1  class QTable:
2      def __init__(self, num_states, num_actions, gamma=0.99):
3          self.gamma = gamma
4          self.Q = np.zeros(shape=(num_states, num_actions))
5
6      def update(self, state, action, reward, next_state, alpha):
7          target = reward + self.gamma * np.max(self.Q[next_state]) - self.Q[state, action]
8          self.Q[state, action] += alpha * target
9
10     def act(self, state):
11         return np.argmax(self.Q[state])
```

store Q-function as a table

Q-learning update!

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t(s_t, a_t) \underbrace{\left(r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)}_{\text{error}}$$

error

greedy action $a_t = \arg \max_a Q(s_t, a)$

Q-learning - Implementation

Complete Outline

```
1 learner = QTable(num_states=env.observation_space.n, num_actions=env.action_space.n, gamma=gamma)
2 rollout_len = 1000000
3 visit_count = np.zeros(shape=(num_states, num_actions)) # save visit counts N(s, a) of all state-action pairs
4 alpha = VisitCountStepsizeSchedule(deg=0.5001)
5 epsilon = LinearExplorationSchedule(rollout_len, final_epsilon=0.4)
6
7 s = env.reset()
8 for t in tqdm(range(rollout_len + 1)):
9     u = np.random.rand()
10    if u < epsilon(t):
11        a = env.action_space.sample()
12    else:
13        a = learner.act(state=s)
14    s_next, r, _, _ = env.step(action=a)
15    n = visit_count[s, a]
16    learner.update(state=s, action=a, reward=r, next_state=s_next, alpha=alpha(n))
17    visit_count[s, a] += 1
18    s = s_next
```

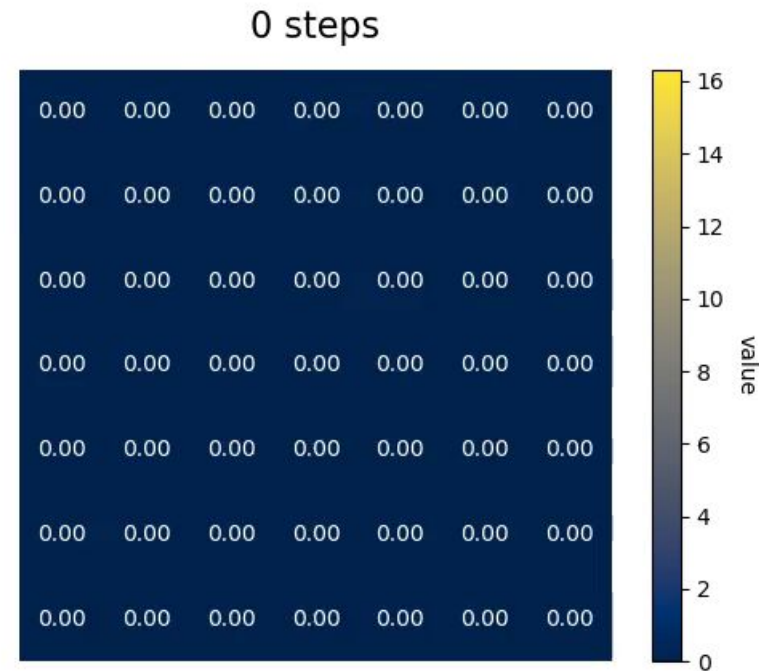
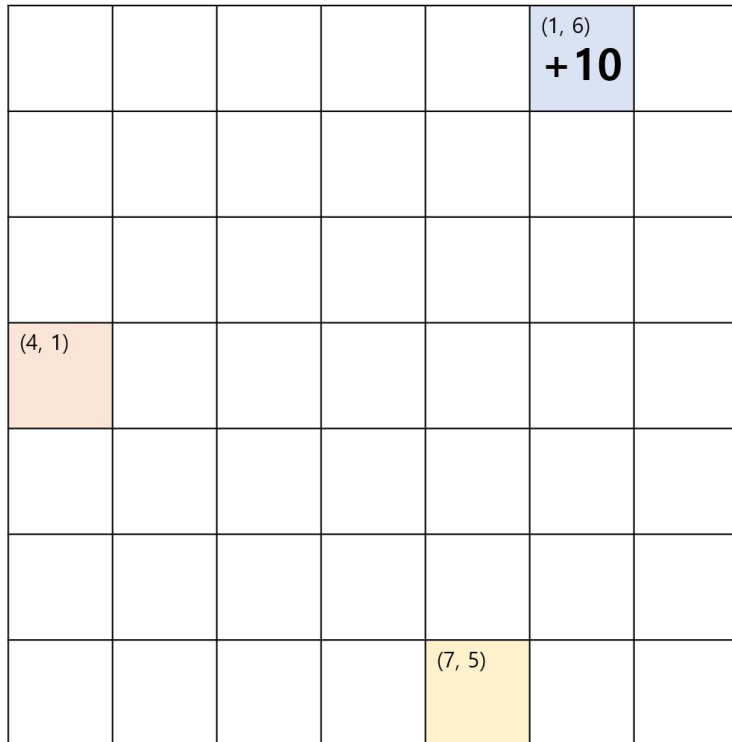
stepsize rule : $\alpha_t(s, a) = \frac{1}{n_t(s, a)^d}$
 $1/2 < d < 1$

exploration strategy : start with large ε , and decrease it.

ε -greedy action selection

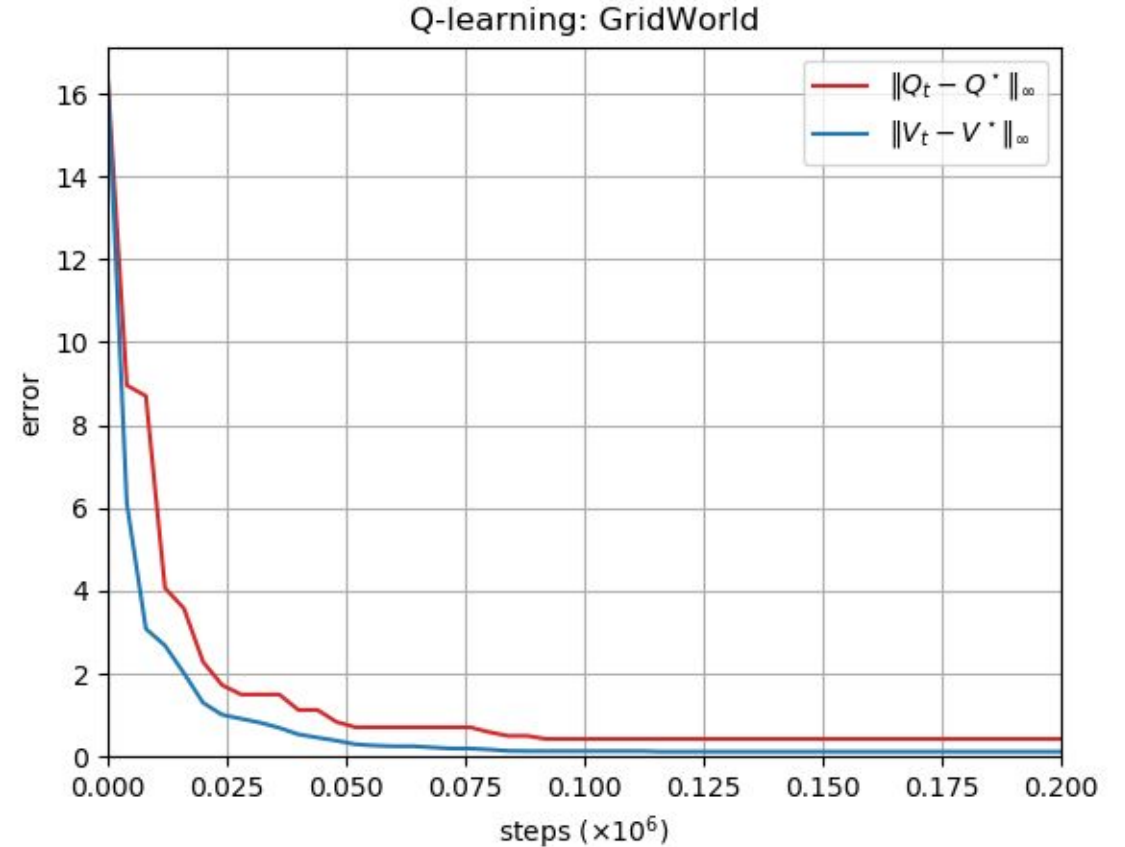
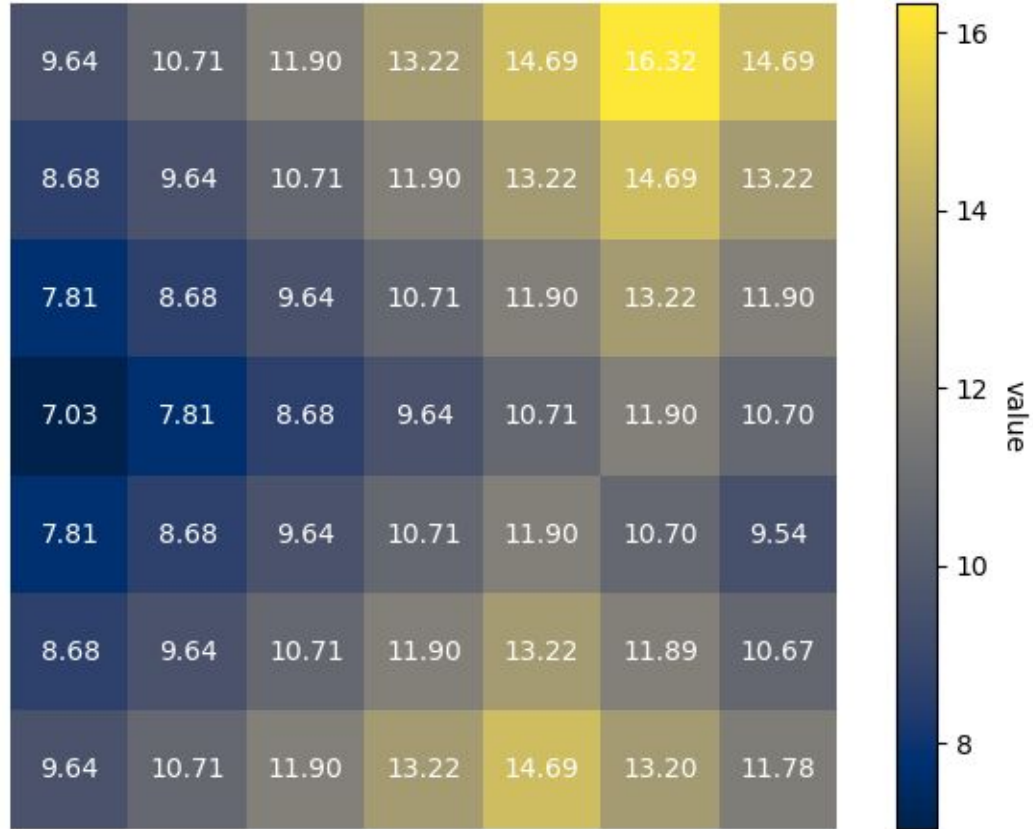
Practice1 - GridWorld

In Day 1, we learned how to compute the value function of GridWorld via **value iteration** when the model is completely **known**.



Practice1 - GridWorld

Q-learning: GridWorld

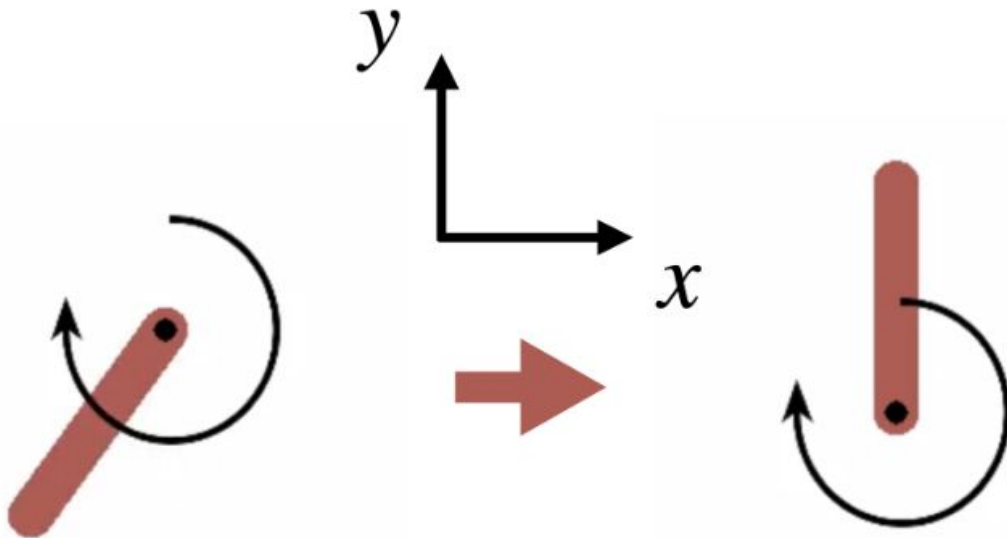


Practice2 – Discrete Pendulum

$$s := (x, y, \dot{\theta}) \quad \begin{cases} x = \cos \theta : x \text{ coordinate of the pendulum tip} \\ y = \sin \theta : y \text{ coordinate of the pendulum tip} \\ \dot{\theta} : \text{Angular velocity} \end{cases}, \quad a := \ddot{\theta},$$

Where the origin is set to the joint of the pendulum, and $-\pi \leq \theta \leq \pi$ as $\theta = 0$ is set to the $+y$ direction. Finally, as control objective is $\theta = \dot{\theta} = 0$, we will give the reward of

$$r := -(\theta^2 + 0.1\dot{\theta}^2 + 0.001\ddot{\theta}^2).$$



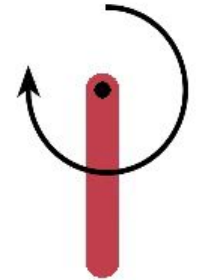
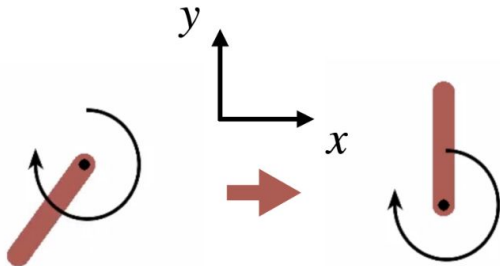
Practice2 – Discrete Pendulum

- discretize pendulum dynamics
 - 2-dim. state $(\theta, \dot{\theta})$
 - 1-dim action τ
- goal : apply a torque τ to a joint for swing-up
- 2460 discretized states & 15 discretized actions

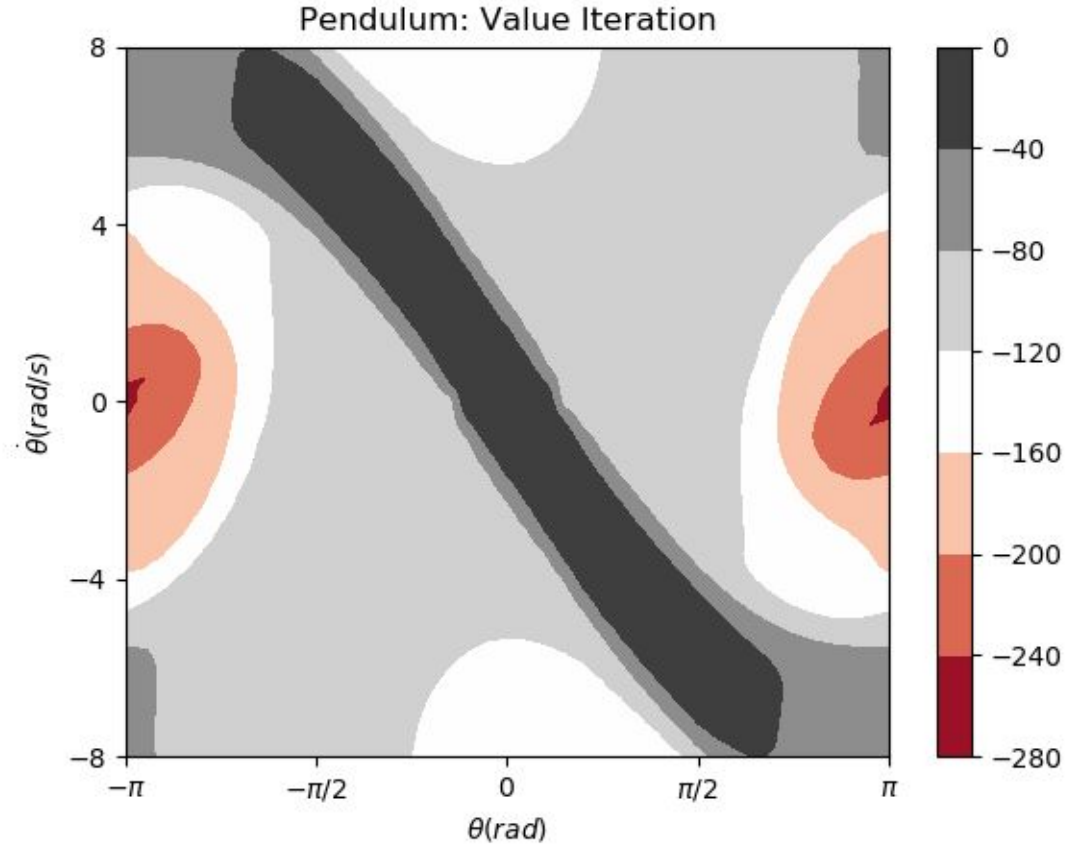
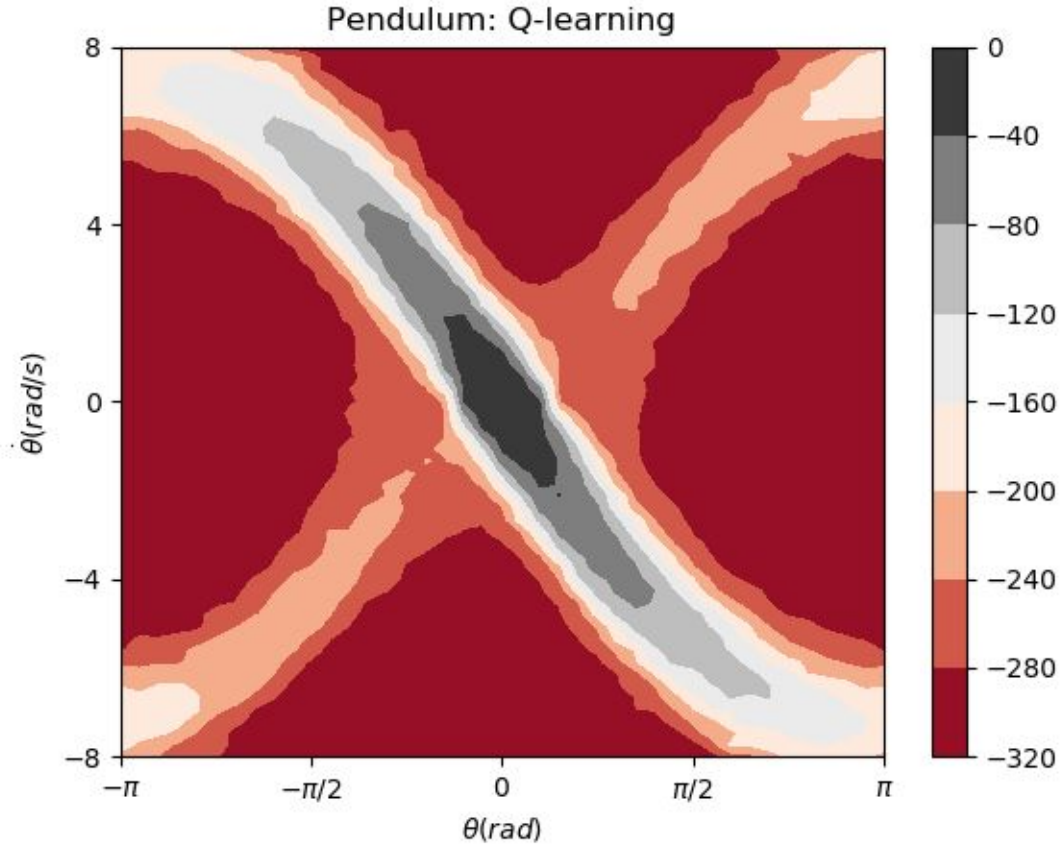
$$s := (x, y, \dot{\theta}) \quad \begin{cases} x = \cos \theta : x \text{ coordinate of the pendulum tip} \\ y = \sin \theta : y \text{ coordinate of the pendulum tip} \\ \dot{\theta} : \text{Angular velocity} \end{cases}, \quad a := \ddot{\theta},$$

Where the origin is set to the joint of the pendulum, and $-\pi \leq \theta \leq \pi$ as $\theta = 0$ is set to the $+y$ direction. Finally, as control objective is $\theta = \dot{\theta} = 0$, we will give the reward of

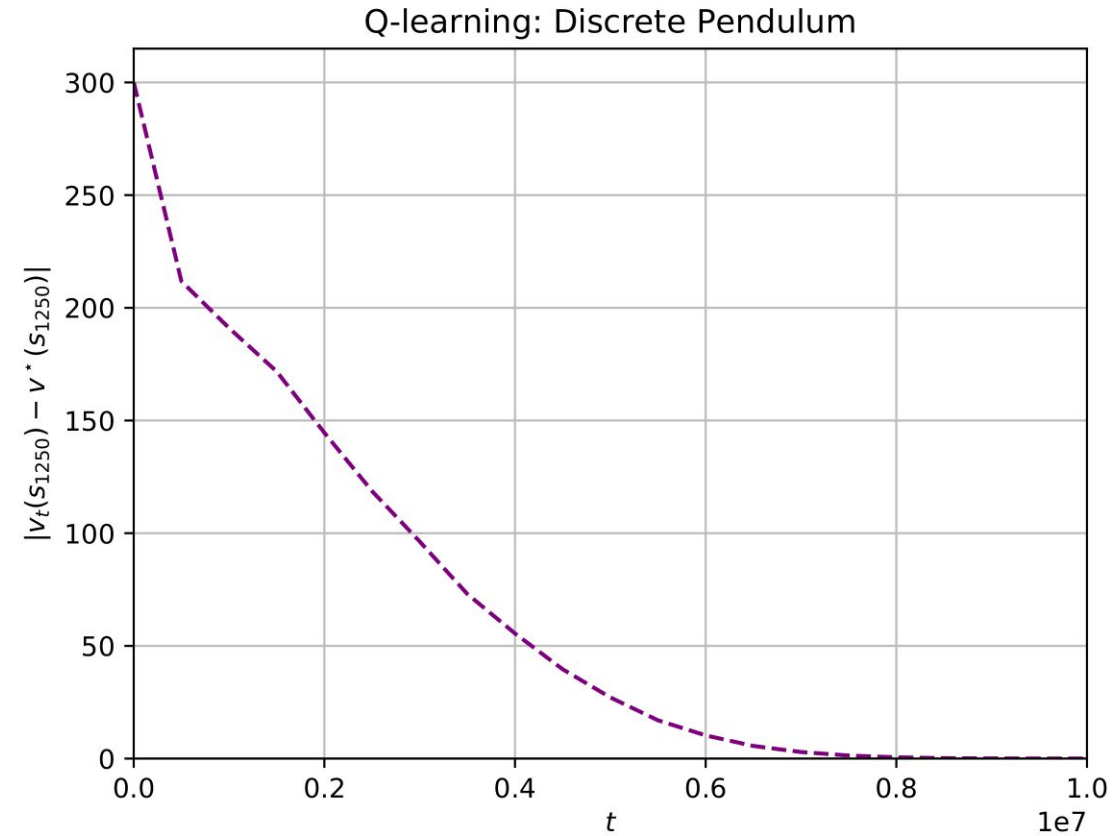
$$r := -(\theta^2 + 0.1\dot{\theta}^2 + 0.001\ddot{\theta}^2).$$



Practice2 – Discrete Pendulum



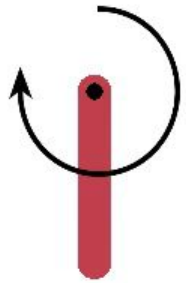
Practice2 – Discrete Pendulum



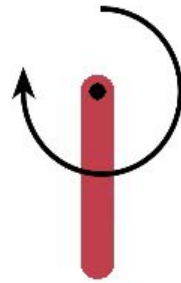
Why Deep Q-Network?

- It seems like Q-learning works well in these examples.
- Why **deep reinforcement learning** then?

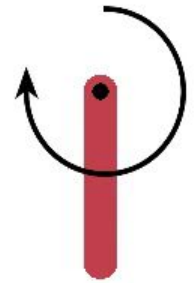
⇒ For most real-world problems, discretization is not a good strategy...



$n = 160, m = 5$ (**coarse**)



$n = 620, m = 10$



$n = 2460, m = 15$ (**fine**)



CORE
Control + Optimization Research Lab