

IERG 5350 Assignment 1: Tabular Reinforcement Learning

2020-2021 Term 1, IERG 5350: Reinforcement Learning. Department of Information Engineering, The Chinese University of Hong Kong. Course Instructor: Professor ZHOU Bolei. Assignment author: PENG Zhenghao, SUN Hao, ZHAN Xiaohang.

Student Name	Student ID
--------------	------------

Wang Wanli	1155160517
------------	------------

Welcome to the assignment 1 of our RL course. The objective of this assignment is for you to understand the classic methods used in tabular reinforcement learning.

This assignment has the following sections:

- Section 1: Warm-up on the RL environment (35 points)
- Section 2: Implementation of model-based family of algorithms: policy iteration and value iteration. (65 points)

You need to go through this self-contained notebook, which contains 21 TODOs in part of the cells and has special `[TODO]` signs. You need to finish all TODOs. Some of them may be easy such as uncommenting a line, some of them may be difficult such as implementing a function. You can find them by searching the `[TODO]` symbol. However, we suggest you to go through the documents step by step, which will give you a better sense of the content.

You are encouraged to add more code on extra cells at the end of the each section to investigate the problems you think interesting. At the end of the file, we left a place for you to optionally write comments (Yes, please give us some either negative or positive rewards so we can keep improving the assignment!).

Please report any code bugs to us via [github issues](#).

Before you get start, remember to follow the instruction at <https://github.com/cuhkrlcourse/ierg5350-assignment> to setup your environment.

Now start running the cells sequentially (by `ctrl + enter` or `shift + enter`) to avoid unnecessary errors by skipping some cells.

Section 1: Warm-up on the RL environment

(35/100 points)

In this section, we will go through the basic concepts of RL environments using OpenAI Gym. Besides, you will get the first sense of the toy environment we will use in the rest of the assignment.

Every Gym environment should contain the following attributes:

1. `env.step(action)` To step the environment by applying `action`. Will return four things: `observation`, `reward`, `done`, `info`, wherein `done` is a boolean value indicating whether this episode is finished. `info` may contain some information the user is interested in, we do not use it.
2. `env.reset()` To reset the environment, back to the initial state. Will return the initial observation.
3. `env.render()` To render the current state of the environment for human-being
4. `env.action_space` The allowed action format. In our case, it is `Discrete(4)` which means the action is an integer in the range [0, 1, 2, 3]. Therefore the `action` for `step(action)` should obey the limit of the action space.
5. `env.observation_space` The observation space.
6. `env.seed(seed)` To set the random seed of the environment. So the result is replicable.

Note that the word episode means the process that an agent interacts with the environment from the initial state to the terminal state. Within one episode, the agent will only receive one `done=True`, when it goes to the terminal state (the agent is dead or the game is over).

We will use "FrozenLake8x8-v0" as our environment. In this environment, the agent controls the movement of a character in a grid world. Some tiles of the grid are walkable, and others lead to the agent falling into the water. Additionally, the movement direction of the agent is uncertain and only partially depends on the chosen direction. The agent is rewarded for finding a walkable path to a goal tile. The meaning of each character:

1. S : starting point, safe
2. F : frozen surface, safe
3. H : hole, fall to your doom
4. G : goal, where the frisbee is located

```
In [1]: from utils import *
import gym
import numpy as np
from collections import deque
```

Section 1.1: Make the environment

You need to know

1. How to make an environment
2. How to set the random seed of environment
3. What is observation space and action space

```
In [2]: env = gym.make('FrozenLake8x8-v0')
env.reset() # [TODO] uncomment this line
env.seed(0) # [TODO] uncomment this line
```

```
print("Current observation space: {}".format(env.observation_space))
print("Current action space: {}".format(env.action_space))
print("0 in action space? {}".format(env.action_space.contains(0)))
print("5 in action space? {}".format(env.action_space.contains(5)))
```

```
Current observation space: Discrete(64)
Current action space: Discrete(4)
0 in action space? True
5 in action space? False
```

Section 1.2: Play the environment with random actions

You need to know

1. How to step the environment
2. How to render the environment

```
In [3]: env.reset()
while True:
    obs, reward, done, info = env.step(env.action_space.sample())# [TODO]
    Uncomment next line
    env.render() # [TODO] Uncomment this line

    print("Current observation: {}\nCurrent reward: {}\n"
          "Whether we are done: {}\ninfo: {}".format(
            obs, reward, done, info
        ))
    wait(sleep=0.5)

    if done == True:
        break
```

```
(Right)
SFFFFFFF
FFFFFFFF
FFFHFFFF
FFFFFHFF
FFFHFFFF
FHHFFHFH
FHFFHFHF
FFFHFFFG
Current observation: 19
Current reward: 0.0
Whether we are done: True
info: {'prob': 0.3333333333333333}
```

Section 1.3: Define the evaluation function to value the random baseline

Now we need to define an evaluation function to evaluate a given policy (a function where the input is observation and the output is action). This is convenient for future evaluation.

As a reminder, you should create a `FrozenLake8x8-v0` environment instance by default, reset it

after each episode (and at the beginning), step the environment, and terminate episode if done.

After implementing the `evaluate` function, run the next cell to check whether you are right.

```
In [4]: def _render_helper(env):
        env.render()
        wait(sleep=0.2)

def evaluate(policy, num_episodes, seed=0, env_name='FrozenLake8x8-v0', render=False):
    """[TODO] You need to implement this function by yourself. It
    evaluate the given policy and return the mean episode reward.
    We use `seed` argument for testing purpose.
    You should pass the tests in the next cell.

    :param policy: a function whose input is an interger (observation)
    :param num_episodes: number of episodes you wish to run
    :param seed: an interger, used for testing.
    :param env_name: the name of the environment
    :param render: a boolean flag. If true, please call _render_helper
    function.
    :return: the averaged episode reward of the given policy.
    """
    env = gym.make(env_name)
    env.seed(seed)

    # Build inner loop to run.
    # For each episode, do not set the limit.
    # Only terminate episode (reset environment) when done = True.
    # The episode reward is the sum of all rewards happen within one episode.

    # Call the helper function `render(env)` to render
    rewards = []
    for i in range(num_episodes):
        obs = env.reset()
        ep_reward = 0
        while True:
            if callable(policy):
                act = policy(obs)
            else:
                act = policy[obs]
            obs, reward, isOver, _ = env.step(act)
            ep_reward += reward
            if render:
                _render_helper(env)
            if isOver:
                break
        rewards.append(ep_reward)
    return np.mean(rewards)

# [TODO] Run next cell to test your implementation!
```

```
In [5]: # Run this cell without modification

# Run this cell to test the correctness of your implementation of `evaluate`.
```

```

LEFT = 0
DOWN = 1
RIGHT = 2
UP = 3

def expert(obs):
    """Go down if agent at the right edge, otherwise go right."""
    return DOWN if (obs + 1) % 8 == 0 else RIGHT

def assert_equal(seed, value, env_name):
    ret = evaluate(expert, 1000, seed, env_name=env_name)
    assert ret == value, \
        "When evaluate on seed {}, 1000 episodes, in {} environment, the " \
        "averaged reward should be {}. But you get {}." \
        "".format(seed, env_name, value, ret)

assert_equal(0, 0.065, 'FrozenLake8x8-v0')
assert_equal(1, 0.059, 'FrozenLake8x8-v0')
assert_equal(2, 0.055, 'FrozenLake8x8-v0')

assert_equal(0, 0.026, 'FrozenLake-v0')
assert_equal(1, 0.034, 'FrozenLake-v0')
assert_equal(2, 0.028, 'FrozenLake-v0')

print("Test Passed!")
print("\nAs a baseline, the mean episode reward of a hand-craft "
      "agent is: ", evaluate(expert, 1000))

```

Test Passed!

As a baseline, the mean episode reward of a hand-craft agent is: 0.065

Congraduation! You have finished section 1 (if and only if not error happen at the above codes).

If you want to do more investigation, feel free to open new cells via `Esc + B` after the next cells and write codes in it, so that you can reuse some result in this notebook. Remember to write sufficient comments and documents to let others know what you are doing.

```

In [6]: ##Details about environment:https://gym.openai.com/docs/
#1: The process gets started by calling reset(), which returns an initial o
bservation.
# observation = env.reset()
#2: The environment's step function returns exactly what we need. In fact,
step
#returns four values. These are:observation (object): an environment-specif
ic object
#representing your observation of the environment. For example, pixel data
from a
#camera, joint angles and joint velocities of a robot, or the board state
in a board
#game.reward (float): amount of reward achieved by the previous action. The
scale
#varies between environments, but the goal is always to increase your total
reward.
#done (boolean): whether it's time to reset the environment again. Most (bu
t not

```

```
#all) tasks are divided up into well-defined episodes, and done being True
indicates
#the episode has terminated. (For example, perhaps the pole tipped too far,
or you
#lost your last life.)info (dict): diagnostic information useful for debugg
ing.
#It can sometimes be useful for learning (for example, it might contain the
raw
#probabilities behind the environment's last state change). However, offici
al
#evaluations of your agent are not allowed to use this for learning.Details
about
#reinforcement learning:https://towardsdatascience.com/reinforcement-learnin
g-
#temporal-difference-sarsa-q-learning-expected-sarsa-on-python-9fecfda746
7e.
##Details about reinforcement learning:https://towardsdatascience.com/reinf
orcement-
#learning-temporal-difference-sarsa-q-learning-expected-sarsa-on-python-9
fecfda7467e
```

Section 2: Model-based Tabular RL

(65/100 points)

We have learned how to use the Gym environment to run an episode, as well as how to interact between the agent (policy) and environment via `env.step(action)` to collect observation, reward, done, and possible extra information.

Now we need to build the basic tabular RL algorithm to solve this environment. Note that compared to the model-free methods in the Sec.3, the algorithms in this section needs to access the internal information of the environment, namely the transition dynamics. In our case, given a state and an action, we need to know which state current environment would jump to, and the probability of this happens, and the reward if the transition happens. You will see that we provide you a helper function `trainer._get_transitions(state, action)` that takes state and action as input and return you a list of possible transitions.

You will use a class to represent a Trainer, which seems to be over-complex for tabular RL. But we will use the same framework in the future assignments, or even in your future research. So it would be helpful for you to get familiar with how to implement an RL algorithm in a class-oriented programming style, as a first step toward the implementation of state of the art RL algorithm in the future.

```
In [7]: # Run this cell without modification

class TabularRLTrainerAbstract:
    """This is the abstract class for tabular RL trainer. We will inherent
the specify
algorithm's trainer from this abstract class, so that we can reuse the
codes like
```

```

    getting the dynamic of the environment (self._get_transitions()) or rendering the
    learned policy (self.render())."""

def __init__(self, env_name='FrozenLake8x8-v0', model_based=True):
    self.env_name = env_name
    self.env = gym.make(self.env_name)
    self.action_dim = self.env.action_space.n
    self.obs_dim = self.env.observation_space.n
    self.model_based = model_based

def _get_transitions(self, state, act):
    """Query the environment to get the transition probability,
    reward, the next state, and done given a pair of state and action.
    We implement this function for you. But you need to know the
    return format of this function.
    """
    self._check_env_name()
    assert self.model_based, "You should not use _get_transitions in "

    "model-free algorithm!"

    # call the internal attribute of the environments.
    # `transitions` is a list contain all possible next states and the
    # probability, reward, and termination indicator corresponding to i
    t
    transitions = self.env.env.P[state][act]

    # Given a certain state and action pair, it is possible
    # to find there exist multiple transitions, since the
    # environment is not deterministic.
    # You need to know the return format of this function: a list of di
    cts

    ret = []
    for prob, next_state, reward, done in transitions:
        ret.append({
            "prob": prob,
            "next_state": next_state,
            "reward": reward,
            "done": done
        })
    return ret

def _check_env_name(self):
    assert self.env_name.startswith('FrozenLake')

def print_table(self):
    """print beautiful table, only work for FrozenLake8X8-v0 env. We
    write this function for you."""
    self._check_env_name()
    print_table(self.table)

def train(self):
    """Conduct one iteration of learning."""
    raise NotImplementedError("You need to override the "
                               "Trainer.train() function.")

```

```

def evaluate(self):
    """Use the function you write to evaluate current policy.
    Return the mean episode reward of 1000 episodes when seed=0."""
    result = evaluate(self.policy, 1000, env_name=self.env_name)
    return result

def render(self):
    """Reuse your evaluate function, render current policy
    for one episode when seed=0"""
    evaluate(self.policy, 1, render=True, env_name=self.env_name)

```

Section 2.1: Policy Iteration

Recall the idea of policy iteration:

1. Update the state value function, given all possible transitions in the environment.
2. Find the best policy that makes the most out of the current state value function.
3. If the best policy is identical to the previous one then stop the training. Otherwise, go to step 1.

In step 1, the way to update the state value function is by

$$v_{k+1}(s) = E_{\pi}[r(s, a) + \gamma v_k(s)]$$

wherein the π is given by current policy, s is next state, r is the reward, $v_k(s)$ is the next state value given by the old (not updated yet) value function. The expectation is computed among all possible transitions (given a state and action pair, it is possible to have many different next states, since the environment is not deterministic).

In step 2, the best policy is to take the action with maximum expected return given a state:

$$\pi(s) = \arg\max_a E_{\pi}[r(s, a) + \gamma v_k(s)]$$

Policy iteration algorithm has an outer loop (update policy, step 1 to 3) and an inner loop (fit the value function, within step 1). In each outer loop, we call once `trainer.train()`, where we call `trainer.update_value_function()` once to update the value function (the state value table). After that we call `trainer.update_policy()` to update the current policy. `trainer` object has a `trainer.policy` attribute, which is a function that takes observation as input and returns an action.

You should implement the trainer following the framework we already wrote for you. Please carefully go through the codes and finish all `TODO` in it.

```

In [8]: # Solve the TODOs and remove `pass`

class PolicyIterationTrainer(TabularRLTrainerAbstract):
    def __init__(self, gamma=1.0, eps=1e-10, env_name='FrozenLake8x8-v0'):
        super(PolicyIterationTrainer, self).__init__(env_name)
        self.gamma = gamma
        self.eps = eps
        self.policy = np.random.choice(self.action_dim, size=(self.obs_dim))

    # test your random policy
    test_random_policy(self.policy, self.env)

    def train(self):
        self.table = np.zeros((self.obs_dim,))

```



```

        self.update_value_function()
        self.update_policy()

    def update_value_function(self):
        count = 0 # count the steps of value updates
        while True:
            old_table = self.table.copy()
            for state in range(self.obs_dim):
                act = self.policy[state]
                transition_list = self._get_transitions(state, act)
                #####
                state_value = []
                for transition in transition_list:
                    prob = transition['prob']
                    reward = transition['reward']
                    next_state = transition['next_state']
                    done = transition['done']
                    state_value.append(prob * (reward + self.gamma*old_table[next_state]))
                self.table[state] = np.sum(state_value)

            if np.sum(np.fabs(old_table - self.table)) <= self.eps:
                break
            count += 1
            if count % 200 == 0:
                # disable this part if you think debug message annoying.
                print("[DEBUG]\tUpdated values for {} steps. "
                      "Difference between new and old table is: {}".format(
                          count, np.sum(np.abs(old_table - self.table))
                      ))
            if count > 4000:
                print("[HINT] Are you sure your codes is OK? It shouldn't be "
                      "so hard to update the value function. You already "
                      "use {} steps to update value function within "
                      "single iteration.".format(count))
            if count > 6000:
                raise ValueError("Clearly your code has problem. Check it!")
        )

    def update_policy(self):
        policy_table = np.zeros([self.obs_dim, ], dtype=np.int)
        for state in range(self.obs_dim):
            state_action_values = []
            for action in range(self.action_dim):
                all_outcomes = []
                transition_list = self._get_transitions(state, action)
                for transition in transition_list:
                    prob = transition['prob']
                    reward = transition['reward']
                    next_state = transition['next_state']
                    done = transition['done']
                    all_outcomes.append(prob * (reward + self.gamma * self.table[next_state]))
                state_action_values.append(np.sum(all_outcomes))
            policy_table[state] = np.argmax(state_action_values)
        self.policy = policy_table

```

Now we have built the Trainer class for policy iteration algorithm. In the following few cells, we will train the agent to solve the problem and evaluate its performance.

```
In [9]: # Solve the TODOs and remove `pass`

# Managing configurations of your experiments is important for your research.
default_pi_config = dict(
    max_iteration=1000,
    evaluate_interval=1,
    gamma=1.0,
    eps=1e-10
)

def policy_iteration(train_config=None):
    config = default_pi_config.copy()
    if train_config is not None:
        config.update(train_config)

    trainer = PolicyIterationTrainer(gamma=config['gamma'], eps=config['eps'])

    old_policy_result = {
        obs: -1 for obs in range(trainer.obs_dim)
    }

    for i in range(config['max_iteration']):
        # train the agent
        trainer.train() # [TODO] please uncomment this line

        # [TODO] compare the new policy with old policy to check whether
        # should we stop. If new and old policy have same output given any
        # observation, then we consider the algorithm is converged and
        # should be stopped.
        if (np.all(old_policy_result == trainer.policy)):
            print("We found policy is not changed anymore at "
                  "iteration {}. Current mean episode reward "
                  "is {}. Stop training.".format(i, trainer.evaluate()))
            break
        old_policy_result = trainer.policy

    # evaluate the result
    if i % config['evaluate_interval'] == 0:
        print(
            "[INFO]\tIn {} iteration, current mean episode reward is {}".
            format(i, trainer.evaluate()))

        if i > 20:
            print("You sure your codes is OK? It shouldn't take so many
                  "
                  "({}) iterations to train a policy iteration "
                  "agent.".format(i))
```

```

    assert trainer.evaluate() > 0.8, \
        "We expect to get the mean episode reward greater than 0.8. " \
        "But you get: {}".format(trainer.evaluate
    ())

    return trainer

```

In [10]: *# Run this cell without modification*

```

# It may be confusing to call a trainer agent. But that's what we normally
do.
pi_agent = policy_iteration()

[INFO] In 0 iteration, current mean episode reward is 0.0.
[INFO] In 1 iteration, current mean episode reward is 0.0.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.0065994699709048255
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.00035060475657588963
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 1.8626298154023058e-05
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 9.89544427232314e-07
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 5.25707344323223e-08
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 2.792882813906772e-09
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 1.483753120368192e-10
[INFO] In 2 iteration, current mean episode reward is 0.0.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.01064319560233583
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.0005654325301605277
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 3.0039281256930006e-05
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 1.595872841673139e-06
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 8.478265914498018e-08
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 4.5041800532386045e-09
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 2.392900703973577e-10
[INFO] In 3 iteration, current mean episode reward is 0.0.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.016743973761202728
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.0008895436862262698
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 4.725807511118951e-05
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 2.5106419144910186e-06
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 1.3338086262901228e-07
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 7.086017848878523e-09

```

[DEBUG] Updated values for 1400 steps. Difference between new and old table is: 3.764535905226296e-10

[INFO] In 4 iteration, current mean episode reward is 0.0.

[DEBUG] Updated values for 200 steps. Difference between new and old table is: 0.01871376846322951

[DEBUG] Updated values for 400 steps. Difference between new and old table is: 0.0009941913920536866

[DEBUG] Updated values for 600 steps. Difference between new and old table is: 5.2817609981757954e-05

[DEBUG] Updated values for 800 steps. Difference between new and old table is: 2.8059988719554563e-06

[DEBUG] Updated values for 1000 steps. Difference between new and old table is: 1.4907205536998647e-07

[DEBUG] Updated values for 1200 steps. Difference between new and old table is: 7.919631178285313e-09

[DEBUG] Updated values for 1400 steps. Difference between new and old table is: 4.2073997237601635e-10

[INFO] In 5 iteration, current mean episode reward is 0.0.

[DEBUG] Updated values for 200 steps. Difference between new and old table is: 0.021815619452000615

[DEBUG] Updated values for 400 steps. Difference between new and old table is: 0.0011589809460460909

[DEBUG] Updated values for 600 steps. Difference between new and old table is: 6.15722526603503e-05

[DEBUG] Updated values for 800 steps. Difference between new and old table is: 3.271099763072806e-06

[DEBUG] Updated values for 1000 steps. Difference between new and old table is: 1.7378109798491664e-07

[DEBUG] Updated values for 1200 steps. Difference between new and old table is: 9.232328391367561e-09

[DEBUG] Updated values for 1400 steps. Difference between new and old table is: 4.904791676973996e-10

[INFO] In 6 iteration, current mean episode reward is 0.615.

[DEBUG] Updated values for 200 steps. Difference between new and old table is: 0.03010397131874412

[DEBUG] Updated values for 400 steps. Difference between new and old table is: 0.0015993389137119486

[DEBUG] Updated values for 600 steps. Difference between new and old table is: 8.496679930397252e-05

[DEBUG] Updated values for 800 steps. Difference between new and old table is: 4.513963109936681e-06

[DEBUG] Updated values for 1000 steps. Difference between new and old table is: 2.398097030467672e-07

[DEBUG] Updated values for 1200 steps. Difference between new and old table is: 1.2740177232628547e-08

[DEBUG] Updated values for 1400 steps. Difference between new and old table is: 6.768373674026495e-10

[INFO] In 7 iteration, current mean episode reward is 0.85.

[DEBUG] Updated values for 200 steps. Difference between new and old table is: 0.04768196238494006

[DEBUG] Updated values for 400 steps. Difference between new and old table is: 0.002536915521952815

[DEBUG] Updated values for 600 steps. Difference between new and old table is: 0.00013477678100466683

[DEBUG] Updated values for 800 steps. Difference between new and old table is: 7.160178124993033e-06

[DEBUG] Updated values for 1000 steps. Difference between new and old table

```
is: 3.803930496842911e-07
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 2.020883657638528e-08
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 1.0736192351812335e-09
[INFO] In 8 iteration, current mean episode reward is 0.858.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.05999088346593476
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.003210333557833303
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 0.00017055529003209835
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 9.06095458003131e-06
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 4.813740787368292e-07
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 2.5573574785608955e-08
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 1.3586263919451724e-09
[INFO] In 9 iteration, current mean episode reward is 0.867.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.0603182991727987
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.0032462018258551523
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 0.00017247142491326473
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 9.162755269701006e-06
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 4.867823626408807e-07
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 2.5860896021856483e-08
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 1.3738913895444682e-09
[INFO] In 10 iteration, current mean episode reward is 0.867.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.06581673541224373
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.0037870765740912193
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 0.00020252141149340597
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 1.076528471952165e-05
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 5.719465970277504e-07
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 3.0385475899286085e-08
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 1.614265332516851e-09
[INFO] In 11 iteration, current mean episode reward is 0.867.
[DEBUG] Updated values for 200 steps. Difference between new and old table
is: 0.07154226746958085
[DEBUG] Updated values for 400 steps. Difference between new and old table
is: 0.004937646720823655
[DEBUG] Updated values for 600 steps. Difference between new and old table
is: 0.00028502578670354384
```

```
[DEBUG] Updated values for 800 steps. Difference between new and old table
is: 1.5596292745820306e-05
[DEBUG] Updated values for 1000 steps. Difference between new and old table
is: 8.376470535947922e-07
[DEBUG] Updated values for 1200 steps. Difference between new and old table
is: 4.468243278155093e-08
[DEBUG] Updated values for 1400 steps. Difference between new and old table
is: 2.3774344715121742e-09
[DEBUG] Updated values for 1600 steps. Difference between new and old table
is: 1.2637649260405226e-10
We found policy is not changed anymore at itertaion 12. Current mean episod
e reward is 0.867. Stop training.
```

```
In [11]: # Run this cell without modification

print("Your policy iteration agent achieve {} mean episode reward. The opti
mal score "
      "should be almost {}".format(pi_agent.evaluate(), 0.86))

Your policy iteration agent achieve 0.867 mean episode reward. The optimal
score should be almost 0.86.
```

```
In [12]: # Run this cell without modification

pi_agent.render()

(Right)
SFFFFFFF
FFFFFFF
FFFHFFF
FFFFFHFF
FFFHFFF
FHHFFHF
FHHFFHF
FFFHFFFG
```

```
In [13]: # Run this cell without modification

pi_agent.print_table()

+-----+-----+-----State Value Mapping-----+-----+-----+
|      | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-----+-----+-----+-----+-----+-----+-----+-----+
| 0 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
|      |      |      |      |      |      |      |      |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 1 | 1.000|1.000|1.000|1.000|1.000|1.000|1.000|1.000|
|      |      |      |      |      |      |      |      |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 2 | 1.000|0.978|0.926|0.000|0.857|0.946|0.982|1.000|
|      |      |      |      |      |      |      |      |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 3 | 1.000|0.935|0.801|0.475|0.624|0.000|0.945|1.000|
|      |      |      |      |      |      |      |      |      |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 4 | 1.000|0.826|0.542|0.000|0.539|0.611|0.852|1.000|
|      |      |      |      |      |      |      |      |      |
```

5	1.000	0.000	0.000	0.168	0.383	0.442	0.000	1.000
6	1.000	0.000	0.195	0.121	0.000	0.332	0.000	1.000
7	1.000	0.732	0.463	0.000	0.277	0.555	0.777	0.000

Congratulations! You have successfully implemented the policy iteration trainer (if and only if no error happens at the above cells). Few further problems for you to investigate:

- 1. What is the impact of the discount factor gamma?
- 2. What is the impact of the value function convergence criterion epsilon?

If you are interested in doing more investigation (not limited to these two), feel free to open new cells at the end of this notebook and left a clear trace of your thinking and coding, which leads to extra credit if you do a good job. It's an optional job, and you can ignore it.

Now let's continue our journey!

Section 2.2: Value Iteration

Recall the idea of value iteration. We update the state value:

$$V_{k+1}(s) = \max_a E_{s'} [r(s, a) + \gamma V_k(s')]$$

wherein the s' is next state, r is the reward, $V_k(s')$ is the next state value given by the old (not updated yet) value function. The expectation is computed among all possible transitions (given a state and action pair, it is possible to have many different next states, since the environment is not deterministic).

The value iteration algorithm does not require an inner loop. It computes the expected return of all possible actions at a given state and uses the maximum of them as the state value. You can imagine it "pretends" we already have the optimal policy and run policy iteration based on it. Therefore we do not need to maintain a policy object in a trainer. We only need to retrieve the optimal policy using the same rule as policy iteration, given current value function.

You should implement the trainer following the framework we already wrote for you. Please carefully go through the code and finish all `TODO` in it.

```
In [14]: # Solve the TODOs and remove `pass`

class ValueIterationTrainer(PolicyIterationTrainer):
    """Note that we inherit Policy Iteration Trainer, to reuse the
    code of update_policy(). It's same since it get optimal policy from
    current state-value table (self.table).
    """
```

```

def __init__(self, gamma=1.0, env_name='FrozenLake8x8-v0'):
    super(ValueIterationTrainer, self).__init__(gamma, None, env_name)
    self.table = np.zeros((self.obs_dim,))
def train(self):
    """Conduct one iteration of learning."""
    # [TODO] value function may be need to be reset to zeros.
    # if you think it should, than do it. If not, then move on.
    # In value iteration, we do not explicit require a
    # policy instance to run. We update value function
    # directly based on the transitions. Therefore, we
    # don't need to run self.update_policy() in each step.
    self.update_value_function()

def update_value_function(self):
    old_table = self.table.copy()
    for state in range(self.obs_dim):
        Q = []
        state_value = 0
        for action in range(self.action_dim):
            state_valueinher = []
            for transition in self._get_transitions(state, action):
                prob = transition['prob']
                reward = transition['reward']
                next_state = transition['next_state']
                done = transition['done']
                state_valueinher.append(prob * (reward + self.gamma * o
ld_table[next_state]))
            Q.append(np.sum(state_valueinher))
        self.table[state] = np.max(Q)

    # Till now the one step value update is finished.
    # You can see that we do not use a inner loop to update
    # the value function like what we did in policy iteration.
    # This is because to compute the state value, which is
    # a expectation among all possible action given by a
    # specified policy, we **pretend** already own the optimal
    # policy (the max operation).

def evaluate(self):
    """Since in value itertaion we do not maintain a policy function,
    so we need to retrieve it when we need it."""
    self.update_policy()
    return super().evaluate()

def render(self):
    """Since in value itertaion we do not maintain a policy function,
    so we need to retrieve it when we need it."""
    self.update_policy()
    return super().render()

```

```

In [15]: # Solve the TODOs and remove `pass`

# Managing configurations of your experiments is important for your research.
default_vi_config = dict(
    max_iteration=10000,

```



```

        evaluate_interval=100, # don't need to update policy each iteration
        gamma=1.0,
        eps=1e-10
    )

def value_iteration(train_config=None):
    config = default_vi_config.copy()
    if train_config is not None:
        config.update(train_config)
    # [TODO] initialize Value Iteration Trainer. Remember to pass
    # config['gamma'] to it.
    trainer = ValueIterationTrainer()
    for i in range(config['max_iteration']):
        old_table = trainer.table.copy()
        # train the agent
        trainer.train() # [TODO] please uncomment this line
        # evaluate the result
        if i % config['evaluate_interval'] == 0:
            print("[INFO]\tIn {} iteration, current "
                  "mean episode reward is {}".format(
                      i, trainer.evaluate()
                  ))

            # [TODO] compare the new policy with old policy to check should
            # we stop.
            # [HINT] If new and old policy have same output given any
            # observation, then we consider the algorithm is converged and
            # should be stopped.

            if(np.sum(np.fabs(old_table - trainer.table)) <= default_vi_con
fig['eps']):
                print("We found policy is not changed anymore at "
                      "iteration {}. Current mean episode reward "
                      "is {}. Stop training.".format(i, trainer.evaluate()))
            )

            break

        if i > 3000:
            print("You sure your codes is OK? It shouldn't take so many
"
                  "({}) iterations to train a policy iteration "
                  "agent.".format(
                      i))

    assert trainer.evaluate() > 0.8, \
        "We expect to get the mean episode reward greater than 0.8. " \
        "But you get: {}. Please check your codes.".format(trainer.evaluate
    ())

    return trainer

```

In [16]: # Run this cell without modification

```
vi_agent = value_iteration()
```

```
[INFO] In 0 iteration, current mean episode reward is 0.0.
```

```
[INFO] In 100 iteration, current mean episode reward is 0.892.
[INFO] In 200 iteration, current mean episode reward is 0.867.
[INFO] In 300 iteration, current mean episode reward is 0.867.
[INFO] In 400 iteration, current mean episode reward is 0.867.
[INFO] In 500 iteration, current mean episode reward is 0.867.
[INFO] In 600 iteration, current mean episode reward is 0.867.
[INFO] In 700 iteration, current mean episode reward is 0.867.
[INFO] In 800 iteration, current mean episode reward is 0.867.
[INFO] In 900 iteration, current mean episode reward is 0.867.
[INFO] In 1000 iteration, current mean episode reward is 0.867.
[INFO] In 1100 iteration, current mean episode reward is 0.867.
[INFO] In 1200 iteration, current mean episode reward is 0.867.
[INFO] In 1300 iteration, current mean episode reward is 0.867.
[INFO] In 1400 iteration, current mean episode reward is 0.867.
[INFO] In 1500 iteration, current mean episode reward is 0.867.
[INFO] In 1600 iteration, current mean episode reward is 0.867.
We found policy is not changed anymore at itertaion 1600. Current mean epis
ode reward is 0.867. Stop training.
```

In [17]:

```
# Run this cell without modification

print("Your value iteration agent achieve {} mean episode reward. The opti
mal score "
      "should be almost {}".format(vi_agent.evaluate(), 0.86))
```

Your value iteration agent achieve 0.867 mean episode reward. The optimal s
core should be almost 0.86.

In [18]:

```
# Run this cell without modification

vi_agent.render()
```

(Right)

```
SFFFFFFF
FFFFFFF
FFFHFFFF
FFFFFHFF
FFFHFFFF
FHHFFHF
FHFFHFHF
FFFHFFFG
```

In [19]:

```
# Run this cell without modification

vi_agent.print_table()
```

	0	1	2	3	4	5	6	7
0	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
1	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
2	1.000	0.978	0.926	0.000	0.857	0.946	0.982	1.000

3	1.000	0.935	0.801	0.475	0.624	0.000	0.945	1.000
4	1.000	0.826	0.542	0.000	0.539	0.611	0.852	1.000
5	1.000	0.000	0.000	0.168	0.383	0.442	0.000	1.000
6	1.000	0.000	0.195	0.121	0.000	0.332	0.000	1.000
7	1.000	0.732	0.463	0.000	0.277	0.555	0.777	0.000

Congratulation! You have successfully implemented the value iteration trainer (if and only if no error happens at the above cells). Few further problems for you to investigate:

- 1. Do you see that some iteration during training yields better rewards than the final one? Why does that happen?
- 2. What is the impact of the discount factor gamma?
- 3. What is the impact of the value function convergence criterion epsilon?

If you are interested in doing more investigation (not limited to these two), feel free to open new cells at the end of this notebook and left a clear trace of your thinking and coding, which leads to extra credit if you do a good job. It's an optional job, and you can ignore it.

Now let's continue our journey!

Section 2.3: Compare two model-based agents

Now we have two agents: `pi_agent` and `vi_agent`. They are believed to be the optimal policy in this environment. Can you compare the policy of two of them and use a clean and clear description or figures to show your conclusion?

```
In [20]: #1: Policy iteration includes policy evaluation and policy iprovement. Thes
e two
#parts are repeated iteratively until policy converges. In contrast, value
iteration includes finding optimal value function
# and policy extraction which is only run one time.
#2: The difference between policy evaluation and finding optimal value func
tion
#except using a max operation for finding optimal value function.
#3: In comparison with the policy iteration starting with a random policy,
the value
#iteration starts with a random value function.
#4: In comparion with the policy iteration, the value iteration is much fas
ter.
```

```
In [21]: # You can do more investigation here if you wish. Leave it blank if you do  
n't.
```

Conclusion and Discussion

In this assignment, we learn how to use Gym package, how to use Object Oriented Programming idea to build a basic tabular RL algorithm.

It's OK to leave the following cells empty. In the next markdown cell, you can write whatever you like. Like the suggestion on the course, the confusing problems in the assignments, and so on.

If you want to do more investigation, feel free to open new cells via `Esc + B` after the next cells and write codes in it, so that you can reuse some result in this notebook. Remember to write sufficient comments and documents to let others know what you are doing.

Following the submission instruction in the assignment to submit your assignment to our staff. Thank you!

...