

IERG 5350 Assignment 2: Model-free Tabular RL

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: Implementation of model-free family of algorithms: SARSA, Q-Learning and model-free control. (100 points)

You need to go through this self-contained notebook, which contains dozens of 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.

Section 1: SARSA

(30/100 points)

You have noticed that in Assignment 1 - Section 2, we always use the function

`trainer._get_transitions()` to get the transition dynamics of the environment, while never call `trainer.env.step()` to really interact with the environment. We need to access the internal feature of the environment or have somebody implement `_get_transitions` for us. However,

this is not feasible in many cases, especially in some real-world cases like autonomous driving where the transition dynamics is unknown or does not explicitly exist.

In this section, we will introduce the Model-free family of algorithms that do not require to know the transitions: they only get information from `env.step(action)`, that collect information by interacting with the environment rather than grab the oracle of the transition dynamics of the environment.

We will continue to use the `TabularRLTrainerAbstract` class to implement algorithms, but remember you should not call `trainer._get_transitions()` anymore.

We will use a simpler environment `FrozenLakerNotSlippery-v0` to conduct experiments, which has a `4 x 4` grids and is deterministic. This is because, in a model-free setting, it's extremely hard for a random agent to achieve the goal for the first time. To reduce the time of experiments, we choose to use a simpler environment. In the bonus section, you will have the chance to try model-free RL on `FrozenLake8x8-v0` to see what will happen.

Now go through each section and start your coding!

Recall the idea of SARSA: it's an on-policy TD control method, which has distinct features compared to policy iteration and value iteration:

1. Maintain a state-action pair value function $Q(s_t, a_t) = E \sum_{i=0} \gamma^i r_{t+i}$, namely the Q value.
2. Do not require to know the internal dynamics of the environment.
3. Use an epsilon-greedy policy to balance exploration and exploitation.

In SARSA algorithm, we update the state action value (Q value) via TD error:

$$TD(s_t, a_t) = r(s_t, a_t) + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

where we run the policy to get the next action $a_{t+1} = \text{Policy}(s_{t+1})$. (That's why we call SARSA an on-policy algorithm, it use the current policy to evaluate Q value).

$$Q^{new}(s_t, a_t) = Q(s_t, a_t) + \alpha TD(s_t, a_t)$$

Wherein α is the learning rate, a hyper-parameter provided by the user.

Now go through the codes.

```
In [25]: # Run this cell without modification

# Import some packages that we need to use
from utils import *
import gym
import numpy as np
from collections import deque
```

```
In [26]: # Solve the TODOs and remove `pass`

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.
    """

    # Create environment (according to env_name, we will use env other than
    'FrozenLake8x8-v0')
    env = gym.make(env_name)

    # Seed the environment
    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 episod
    e.

    # Call the helper function `_render_helper(env)` to render
    rewards = []
    for i in range(num_episodes):
        # reset the environment
        obs = env.reset()
        act = policy(obs)

        ep_reward = 0.0
        while True:
            next_obs, reward, done, _ = env.step(act)
            act = policy(next_obs)

            ep_reward += reward

            if render:
                _render_helper(env)
            if done:
                break

        rewards.append(ep_reward)

    return np.mean(rewards)

```

In [27]: # 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 ren
dering 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)

```

In [28]: # Solve the TODOs and remove `pass`

```

class SARSATrainer(TabularRLTrainerAbstract):
    def __init__(self,
                  gamma=1.0,
                  eps=0.1,
                  learning_rate=1.0,
                  max_episode_length=100,
                  env_name='FrozenLake8x8-v0'
                  ):
        super(SARSATrainer, self).__init__(env_name, model_based=False)

        # discount factor
        self.gamma = gamma

        # epsilon-greedy exploration policy parameter
        self.eps = eps

        # maximum steps in single episode
        self.max_episode_length = max_episode_length

        # the learning rate
        self.learning_rate = learning_rate

        # build the Q table
        self.table = np.zeros((self.obs_dim, self.action_dim))

    def policy(self, obs):
        """Implement epsilon-greedy policy

        It is a function that take an integer (state / observation)
        as input and return an interger (action).
        """
        p = np.random.random_sample()
        if p <= self.eps:
            return np.random.randint(self.action_dim)
        else:
            return np.argmax(self.table[obs, :])

    def train(self):
        """Conduct one iteration of learning."""

```

```

        self.eps = 1
        obs = self.env.reset()
        act = self.policy(obs)

        for t in range(self.max_episode_length):
            # Gradually reduce epsilon to the minimum: first exploration, then exploitation
            if self.eps > 0.1:
                self.eps -= 0.01

            next_obs, reward, done, _ = self.env.step(act)
            next_act = self.policy(next_obs)

            td_error = reward + self.gamma * self.table[next_obs, next_act] - self.table[obs, act]

            new_value = self.table[obs, act] + self.learning_rate * td_error

            self.table[obs, act] = new_value

            obs = next_obs
            act = next_act

            if done:
                break

```

Now you have finish the SARSA trainer. To make sure your implementation of epsilon-greedy strategy is correct, please run the next cell.

```

In [29]: # Run this cell without modification

# set eps = 0 to disable exploration.
test_trainer = SARSATrainer(eps=0.0)
test_trainer.table.fill(0)

# set the Q value of (obs 0, act 3) to 100, so that it should be taken by # policy.
test_obs = 0
test_act = test_trainer.action_dim - 1
test_trainer.table[test_obs][test_act] = 100

# assertion
assert test_trainer.policy(test_obs) == test_act, \
    "Your action is wrong! Should be {} but get {}".format(
        test_act, test_trainer.policy(test_obs))

# delete trainer
del test_trainer

# set eps = 0 to disable exploitation.
test_trainer = SARSATrainer(eps=1.0)
test_trainer.table.fill(0)

act_set = set()

```

```

for i in range(100):
    act_set.add(test_trainer.policy(0))

# assertion
assert len(act_set) > 1, ("You sure your uniformly action selection mechanism"
                           " is working? You only take action {} when "
                           "observation is 0, though we run trainer.policy()"
                           "for 100 times.".format(act_set))

# delete trainer
del test_trainer

print("Policy Test passed!")

```

Policy Test passed!

Now run the next cell to see the result. Note that we use the non-slippy version of a small frozen lake environment `FrozenLakeNotSlippery-v0` (this is not a ready Gym environment, see `utils.py` for details). This is because, in the model-free setting, it's extremely hard to access the goal for the first time (you should already know that if you watch the agent randomly acting in Assignment 1 - Section 1).

```

In [30]: # Solve TODO

# Managing configurations of your experiments is important for your research.
default_sarsa_config = dict(
    max_iteration=20000,
    max_episode_length=200,
    learning_rate=0.1,
    evaluate_interval=1000,
    gamma=0.95,
    eps=1,
    env_name='FrozenLakeNotSlippery-v0'
)

def sarsa(train_config=None):
    config = default_sarsa_config.copy()
    if train_config is not None:
        config.update(train_config)

    trainer = SARSATrainer(
        gamma=config['gamma'],
        eps=config['eps'],
        learning_rate=config['learning_rate'],
        max_episode_length=config['max_episode_length'],
        env_name=config['env_name']
    )

    for i in range(1, config['max_iteration'] + 1):
        # train the agent
        trainer.train() # [TODO] please uncomment this line

```

```

# evaluate the result
if i % config['evaluate_interval'] == 0:

    # gradually reduce epsilon to 0 before each evaluation
    trainer.eps = 1e-2 * (config['max_iteration'] - i) / config['e
valuate_interval']

    '''
    --- THIS IS THE KEY! ---
    No matter how well your trainer works, as long as trainer.ep
silon > 0,
    the evaluation function will always randomly select action occa
sionally,
    which extremely hinders the episode rewards.
    With this sentence, the mean episode reward can reach the f
ull score 1.0,
    as a comparison, without this sentence, the reward just fluctua
te around 0.1.
    '''

    print(
        "[INFO]\tIn {} iteration, current mean episode reward is {}".
        format(i, trainer.evaluate()))

# to make sure the evaluate runs normally...
trainer.eps = 0

if trainer.evaluate() < 0.6:
    print("We expect to get the mean episode reward greater than 0.6. "
        "\n        "But you get: {}. Please check your codes.".format(trainer.evaluate
        ()))

return trainer

```

In [31]: # Run this cell without modification

```
sarsa_trainer = sarsa()
```

```

[INFO] In 1000 iteration, current mean episode reward is 0.78.
[INFO] In 2000 iteration, current mean episode reward is 0.804.
[INFO] In 3000 iteration, current mean episode reward is 0.838.
[INFO] In 4000 iteration, current mean episode reward is 0.807.
[INFO] In 5000 iteration, current mean episode reward is 0.852.
[INFO] In 6000 iteration, current mean episode reward is 0.856.
[INFO] In 7000 iteration, current mean episode reward is 0.847.
[INFO] In 8000 iteration, current mean episode reward is 0.854.
[INFO] In 9000 iteration, current mean episode reward is 0.875.
[INFO] In 10000 iteration, current mean episode reward is 0.885.
[INFO] In 11000 iteration, current mean episode reward is 0.893.
[INFO] In 12000 iteration, current mean episode reward is 0.933.
[INFO] In 13000 iteration, current mean episode reward is 0.916.
[INFO] In 14000 iteration, current mean episode reward is 0.937.
[INFO] In 15000 iteration, current mean episode reward is 0.962.
[INFO] In 16000 iteration, current mean episode reward is 0.959.
[INFO] In 17000 iteration, current mean episode reward is 0.978.

```


[INFO] In 18000 iteration, current mean episode reward is 0.825.
[INFO] In 19000 iteration, current mean episode reward is 0.983.
[INFO] In 20000 iteration, current mean episode reward is 1.0.

In [32]: *# Run this cell without modification*

```
sarsa_trainer.print_table()

=== The state value for action 0 ===
+-----+-----+-----+-----+-----+
|       | 0   | 1   | 2   | 3   |
+-----+-----+-----+-----+-----+
| 0      | 0.020|0.019|0.016|0.035|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 1      | 0.040|0.000|0.000|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 2      | 0.055|0.054|0.101|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 3      | 0.000|0.000|0.208|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+

=== The state value for action 1 ===
+-----+-----+-----+-----+-----+
|       | 0   | 1   | 2   | 3   |
+-----+-----+-----+-----+-----+
| 0      | 0.024|0.000|0.053|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 1      | 0.057|0.000|0.244|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 2      | 0.000|0.193|0.600|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 3      | 0.000|0.248|0.578|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+

=== The state value for action 2 ===
+-----+-----+-----+-----+-----+
|       | 0   | 1   | 2   | 3   |
+-----+-----+-----+-----+-----+
| 0      | 0.017|0.027|0.020|0.020|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 1      | 0.000|0.000|0.000|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
| 2      | 0.134|0.251|0.000|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+-----+
```

	3		0.000		0.501		1.000		0.000	
+-----+-----+-----+-----+-----+										

=== The state value for action 3 ===

+-----+-----+-----+-----+-----+										
			0		1		2		3	
	-----+-----+-----+-----+-----+									
	0		0.020		0.017		0.027		0.019	
+-----+-----+-----+-----+-----+										
	1		0.021		0.000		0.027		0.000	
+-----+-----+-----+-----+-----+										
	2		0.028		0.000		0.128		0.000	
+-----+-----+-----+-----+-----+										
	3		0.000		0.125		0.148		0.000	
+-----+-----+-----+-----+-----+										

```
In [33]: # Run this cell without modification

sarsa_trainer.render()

(Right)
SFFF
FHFH
FFFH
HFFG
```

Now you have finished the SARSA algorithm.

Section 2: Q-Learning

(30/100 points)

Q-learning is an off-policy algorithm who differs from SARSA in the computing of TD error. Instead of running policy to get `next_act` `a` and get the TD error by:

$$r + \gamma Q(s', a) - Q(s, a),$$

in Q-learning we compute the TD error via:

$$r + \gamma \max_{a'} Q(s', a') - Q(s, a).$$

The reason we call it "off-policy" is that the policy involves the computing of next-Q value is not the "behavior policy", instead, it is a "virtual policy" that always takes the best action given current Q values.

```
In [34]: # Solve the TODOs and remove `pass`
```

```

class QLearningTrainer(TabularRLTrainerAbstract):
    def __init__(self,
                  gamma=1.0,
                  eps=0.1,
                  learning_rate=1.0,
                  max_episode_length=100,
                  env_name='FrozenLake8x8-v0'
                  ):
        super(QLearningTrainer, self).__init__(env_name, model_based=False)

        self.gamma = gamma
        self.eps = eps
        self.max_episode_length = max_episode_length
        self.learning_rate = learning_rate

        # build the Q table
        self.table = np.zeros((self.obs_dim, self.action_dim))

    def policy(self, obs):
        """Implement epsilon-greedy policy

        It is a function that take an integer (state / observation)
        as input and return an interger (action).
        """
        p = np.random.random_sample()
        if p <= self.eps:
            return np.random.randint(self.action_dim)
        else:
            return np.argmax(self.table[obs,:])

    def train(self):
        """Conduct one iteration of learning."""

        self.eps = 1
        obs = self.env.reset()

        for t in range(self.max_episode_length):
            # Gradually reduce epsilon to the minimum: first exploration, then exploitation
            if self.eps > 0.1:
                self.eps -= 0.01

            act = self.policy(obs)

            next_obs, reward, done, _ = self.env.step(act)

            td_error = reward + self.gamma * np.max(self.table[next_obs,:]) - self.table[obs, act]

            new_value = self.table[obs, act] + self.learning_rate * td_error

            self.table[obs, act] = new_value
            obs = next_obs

            if done:
                break

```

```

In [35]: # Solve the TODO

# Managing configurations of your experiments is important for your research.
default_q_learning_config = dict(
    max_iteration=20000,
    max_episode_length=200,
    learning_rate=0.05,
    evaluate_interval=1000,
    gamma=0.8,
    eps=1,
    env_name='FrozenLakeNotSlippery-v0'
)

def q_learning(train_config=None):
    config = default_q_learning_config.copy()
    if train_config is not None:
        config.update(train_config)

    trainer = QLearningTrainer(
        gamma=config['gamma'],
        eps=config['eps'],
        learning_rate=config['learning_rate'],
        max_episode_length=config['max_episode_length'],
        env_name=config['env_name']
    )

    for i in range(1, config['max_iteration'] + 1):
        # train the agent
        trainer.train() # [TODO] please uncomment this line

        # evaluate the result
        if i % config['evaluate_interval'] == 0:

            # gradually reduce epsilon to 0 before each evaluation
            trainer.eps = 1e-2 * (config['max_iteration'] - i) / config['evaluate_interval']

            '''
            --- THIS IS THE KEY! ---
            No matter how well your trainer works, as long as trainer.epsilon > 0,
            the evaluation function will always randomly select action occasionally,
            which extremely hinders the episode rewards.
            With this sentence, the mean episode reward can reach the full score 1.0,
            as a comparison, without this sentence, the reward just fluctuate around 0.1.
            '''

            print(
                "[INFO]\tIn {} iteration, current mean episode reward is {}".format(
                    i, trainer.evaluate()
                )
            )

```

```

# to make sure the evaluate runs normally...
trainer.eps = 0

if trainer.evaluate() < 0.6:
    print("We expect to get the mean episode reward greater than 0.6. "
\
        "But you get: {}".format(trainer.evaluate
    ()))

return trainer

```

In [36]: *# Run this cell without modification*

```
q_learning_trainer = q_learning()
```

```

[INFO] In 1000 iteration, current mean episode reward is 0.763.
[INFO] In 2000 iteration, current mean episode reward is 0.81.
[INFO] In 3000 iteration, current mean episode reward is 0.807.
[INFO] In 4000 iteration, current mean episode reward is 0.813.
[INFO] In 5000 iteration, current mean episode reward is 0.852.
[INFO] In 6000 iteration, current mean episode reward is 0.843.
[INFO] In 7000 iteration, current mean episode reward is 0.875.
[INFO] In 8000 iteration, current mean episode reward is 0.872.
[INFO] In 9000 iteration, current mean episode reward is 0.868.
[INFO] In 10000 iteration, current mean episode reward is 0.889.
[INFO] In 11000 iteration, current mean episode reward is 0.917.
[INFO] In 12000 iteration, current mean episode reward is 0.918.
[INFO] In 13000 iteration, current mean episode reward is 0.932.
[INFO] In 14000 iteration, current mean episode reward is 0.935.
[INFO] In 15000 iteration, current mean episode reward is 0.935.
[INFO] In 16000 iteration, current mean episode reward is 0.965.
[INFO] In 17000 iteration, current mean episode reward is 0.969.
[INFO] In 18000 iteration, current mean episode reward is 0.974.
[INFO] In 19000 iteration, current mean episode reward is 0.988.
[INFO] In 20000 iteration, current mean episode reward is 1.0.

```

In [37]: *# Run this cell without modification*

```
q_learning_trainer.print_table()
```

```
=== The state value for action 0 ===
```

```

+-----+-----+-----+-----+-----+
|       | 0   | 1   | 2   | 3   |
+-----+-----+-----+-----+
| 0      | 0.262|0.262|0.328|0.410|
|       |      |      |      |      |
+-----+-----+-----+-----+
| 1      | 0.328|0.000|0.000|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+
| 2      | 0.410|0.410|0.512|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+
| 3      | 0.000|0.000|0.640|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+

```

=== The state value for action 1 ===

	0	1	2	3
0	0.328	0.000	0.512	0.000
1	0.410	0.000	0.640	0.000
2	0.000	0.640	0.800	0.000
3	0.000	0.640	0.800	0.000

=== The state value for action 2 ===

	0	1	2	3
0	0.328	0.410	0.328	0.328
1	0.000	0.000	0.000	0.000
2	0.512	0.640	0.000	0.000
3	0.000	0.800	1.000	0.000

=== The state value for action 3 ===

	0	1	2	3
0	0.262	0.328	0.410	0.328
1	0.262	0.000	0.410	0.000
2	0.328	0.000	0.512	0.000
3	0.000	0.512	0.640	0.000

```
In [38]: # Run this cell without modification
```

```
q_learning_trainer.render()
```

```
(Right)
```

```
SFFF
```

```
FHFH
```

```
FFFH
```

```
HFFG
```

Now you have finished Q-Learning algorithm.

Section 3: Monte Carlo Control

(40/100 points)

In sections 1 and 2, we implement the on-policy and off-policy versions of the TD Learning algorithms. In this section, we will play with another branch of the model-free algorithm: Monte Carlo Control. You can refer to the 5.3 Monte Carlo Control section of the textbook "Reinforcement Learning: An Introduction" to learn the details of MC control.

The basic idea of MC control is to compute the Q value (state-action value) directly from an episode, without using TD to fit the Q function. Concretely, we maintain a batch of lists (the total number of lists is `obs_dim * action_dim`), each element of the batch is a list correspondent to a state-action pair. The list is used to store the previously happening "return" of each state action pair.

We will use a dict `self.returns` to store all lists. The keys of the dict are tuples `(obs, act)`: `self.returns[(obs, act)]` is the list to store all returns when `(obs, act)` happens.

The key point of MC Control method is that we take the mean of this list (the mean of all previous returns) as the Q value of this state-action pair.

The "return" here is the discounted return starting from the state-action pair: $\text{Return}(s_t, a_t) = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$.

In short, MC Control method uses a new way to estimate the values of state-action pairs.

```
In [39]: # Solve the TODOs and remove `pass`
```

```
class MCControlTrainer(TabularRLTrainerAbstract):
    def __init__(self,
                  gamma=1.0,
                  eps=0.3,
                  max_episode_length=100,
                  env_name='FrozenLake8x8-v0'):
        super(MCControlTrainer, self).__init__(env_name, model_based=False)

        self.gamma = gamma
        self.eps = eps
        self.max_episode_length = max_episode_length

        # build the dict of lists
```

```

self.returns = {}
for obs in range(self.obs_dim):
    for act in range(self.action_dim):
        self.returns[(obs, act)] = []

# build the Q table
self.table = np.zeros((self.obs_dim, self.action_dim))

def policy(self, obs):
    """Implement epsilon-greedy policy

    It is a function that take an integer (state / observation)
    as input and return an interger (action).
    """
    p = np.random.random_sample()
    if p <= self.eps:
        return np.random.randint(self.action_dim)
    else:
        return np.argmax(self.table[obs,:])

def train(self):
    """Conduct one iteration of learning."""
    observations = []
    actions = []
    rewards = []

    obs = self.env.reset()
    self.eps = 1

    for t in range(self.max_episode_length):

        # Gradually reduce epsilon to the minimum: first exploration, then exploitation
        if self.eps > 0.1:
            self.eps -= 0.01

        act = self.policy(obs)
        next_obs, reward, done, _ = self.env.step(act)

        observations.append(obs)
        actions.append(act)
        rewards.append(reward)

        obs = next_obs

        if done:
            break

    assert len(actions) == len(observations)
    assert len(actions) == len(rewards)

    occured_state_action_pair = set()
    length = len(actions)
    g = 0.0
    for i in reversed(range(length)):
        # if length = 10, then i = 9, 8, ..., 0

```



```

        obs = observations[i]
        act = actions[i]
        reward = rewards[i]

        g = self.gamma * g + reward

        if (obs, act) not in occurred_state_action_pair:
            occurred_state_action_pair.add((obs, act))

            self.returns[(obs, act)].append(g)

            self.table[obs, act] = np.average(self.returns[(obs, act)])
    )

```

```

In [40]: # Run this cell without modification

# Managing configurations of your experiments is important for your research.
default_mc_control_config = dict(
    max_iteration=10000,
    max_episode_length=200,
    evaluate_interval=1000,
    gamma=0.9,
    eps=0.3,
    env_name='FrozenLakeNotSlippery-v0'
)

def mc_control(train_config=None):
    config = default_mc_control_config.copy()
    if train_config is not None:
        config.update(train_config)

    trainer = MCControlTrainer(
        gamma=config['gamma'],
        eps=config['eps'],
        max_episode_length=config['max_episode_length'],
        env_name=config['env_name']
    )

    for i in range(1, config['max_iteration'] + 1):
        # train the agent
        trainer.train()

        # evaluate the result
        if i % config['evaluate_interval'] == 0:

            # gradually reduce epsilon to 0 before each evaluation
            trainer.eps = 1e-2 * (config['max_iteration'] - i) / config['evaluate_interval']

            print(
                "[INFO]\tIn {} iteration, current mean episode reward is {}".format(i, trainer.evaluate())
            )

```

```

# to make sure the evaluate runs normally...
trainer.eps = 0
if trainer.evaluate() < 0.6:
    print("We expect to get the mean episode reward greater than 0.6. "
\
        "But you get: {}".format(trainer.evaluate
    ()))

return trainer

```

In [41]: # Run this cell without modification

```

mc_control_trainer = mc_control()

sarsa_trainer = sarsa()

```

```

[INFO] In 1000 iteration, current mean episode reward is 0.908.
[INFO] In 2000 iteration, current mean episode reward is 0.91.
[INFO] In 3000 iteration, current mean episode reward is 0.927.
[INFO] In 4000 iteration, current mean episode reward is 0.938.
[INFO] In 5000 iteration, current mean episode reward is 0.956.
[INFO] In 6000 iteration, current mean episode reward is 0.967.
[INFO] In 7000 iteration, current mean episode reward is 0.964.
[INFO] In 8000 iteration, current mean episode reward is 0.974.
[INFO] In 9000 iteration, current mean episode reward is 0.993.
[INFO] In 10000 iteration, current mean episode reward is 1.0.
[INFO] In 1000 iteration, current mean episode reward is 0.792.
[INFO] In 2000 iteration, current mean episode reward is 0.804.
[INFO] In 3000 iteration, current mean episode reward is 0.794.
[INFO] In 4000 iteration, current mean episode reward is 0.813.
[INFO] In 5000 iteration, current mean episode reward is 0.837.
[INFO] In 6000 iteration, current mean episode reward is 0.837.
[INFO] In 7000 iteration, current mean episode reward is 0.864.
[INFO] In 8000 iteration, current mean episode reward is 0.882.
[INFO] In 9000 iteration, current mean episode reward is 0.898.
[INFO] In 10000 iteration, current mean episode reward is 0.904.
[INFO] In 11000 iteration, current mean episode reward is 0.904.
[INFO] In 12000 iteration, current mean episode reward is 0.922.
[INFO] In 13000 iteration, current mean episode reward is 0.92.
[INFO] In 14000 iteration, current mean episode reward is 0.932.
[INFO] In 15000 iteration, current mean episode reward is 0.953.
[INFO] In 16000 iteration, current mean episode reward is 0.962.
[INFO] In 17000 iteration, current mean episode reward is 0.964.
[INFO] In 18000 iteration, current mean episode reward is 0.975.
[INFO] In 19000 iteration, current mean episode reward is 0.648.
[INFO] In 20000 iteration, current mean episode reward is 1.0.

```

In [42]: # Run this cell without modification

```

mc_control_trainer.print_table()

```

```

=== The state value for action 0 ===
+-----+-----+-----+-----+-----+
|       | 0   | 1   | 2   | 3   |
|-----+-----+-----+-----+-----+
| 0     | 0.017| 0.017| 0.019| 0.049|
|       |       |       |       |       |

```

1	0.023	0.000	0.000	0.000	
2	0.048	0.045	0.113	0.000	
3	0.000	0.000	0.274	0.000	

=== The state value for action 1 ===

		0	1	2	3
0	0.020	0.000	0.062	0.000	
1	0.046	0.000	0.195	0.000	
2	0.000	0.221	0.560	0.000	
3	0.000	0.230	0.576	0.000	

=== The state value for action 2 ===

		0	1	2	3
0	0.013	0.029	0.023	0.032	
1	0.000	0.000	0.000	0.000	
2	0.117	0.191	0.000	0.000	
3	0.000	0.561	1.000	0.000	

=== The state value for action 3 ===

		0	1	2	3
0	0.016	0.016	0.030	0.026	
1	0.020	0.000	0.037	0.000	

```

+-----+-----+-----+-----+-----+
|  2    |0.025|0.000|0.055|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+
|  3    |0.000|0.119|0.234|0.000|
|       |      |      |      |      |
+-----+-----+-----+-----+

```

In [43]: *# Run this cell without modification*

```
mc_control_trainer.render()
```

```
(Right)
```

```
SFFF
```

```
FHFH
```

```
FFFH
```

```
HFFG
```

Secion 4 Bonus (optional): Tune and train FrozenLake8x8-v0 with Model-free algorithms

You have noticed that we use a simpler environment `FrozenLakeNotSlippery-v0` which has only 16 states and is not stochastic. Can you try to train Model-free families of algorithm using the `FrozenLake8x8-v0` environment? Tune the hyperparameters and compare the results between different algorithms.

Hint: It's not easy to train model-free algorithm in `FrozenLake8x8-v0`. Failure is expected.

In [44]: *# It's ok to leave this cell commented.*

```
new_config = dict(
    env_name="FrozenLake8x8-v0"
)
```

```
# new_mc_control_trainer = mc_control(new_config)
new_q_learning_trainer = q_learning(new_config)
```

```

[INFO] In 1000 iteration, current mean episode reward is 0.0.
[INFO] In 2000 iteration, current mean episode reward is 0.0.
[INFO] In 3000 iteration, current mean episode reward is 0.0.
[INFO] In 4000 iteration, current mean episode reward is 0.0.
[INFO] In 5000 iteration, current mean episode reward is 0.0.
[INFO] In 6000 iteration, current mean episode reward is 0.0.
[INFO] In 7000 iteration, current mean episode reward is 0.033.
[INFO] In 8000 iteration, current mean episode reward is 0.051.
[INFO] In 9000 iteration, current mean episode reward is 0.068.
[INFO] In 10000 iteration, current mean episode reward is 0.172.
[INFO] In 11000 iteration, current mean episode reward is 0.147.
[INFO] In 12000 iteration, current mean episode reward is 0.18.
[INFO] In 13000 iteration, current mean episode reward is 0.4.
[INFO] In 14000 iteration, current mean episode reward is 0.141.
[INFO] In 15000 iteration, current mean episode reward is 0.373.
[INFO] In 16000 iteration, current mean episode reward is 0.253.

```

```
[INFO] In 17000 iteration, current mean episode reward is 0.109.
[INFO] In 18000 iteration, current mean episode reward is 0.13.
[INFO] In 19000 iteration, current mean episode reward is 0.161.
[INFO] In 20000 iteration, current mean episode reward is 0.599.
We expect to get the mean episode reward greater than 0.6. But you get: 0.5
99. Please check your codes.
```

```
In [45]: new_q_learning_trainer.render()
```

```
(Right)
SFFFFFFF
FFFFFFF
FFFHFFFF
FFFFFHFF
FFFHFFFF
FHHFFHF
FHFFHFHF
FFFHFFFG
```

```
In [46]: new_sarsa_trainer = sarsa(new_config)
```

```
[INFO] In 1000 iteration, current mean episode reward is 0.0.
[INFO] In 2000 iteration, current mean episode reward is 0.0.
[INFO] In 3000 iteration, current mean episode reward is 0.0.
[INFO] In 4000 iteration, current mean episode reward is 0.0.
[INFO] In 5000 iteration, current mean episode reward is 0.039.
[INFO] In 6000 iteration, current mean episode reward is 0.129.
[INFO] In 7000 iteration, current mean episode reward is 0.244.
[INFO] In 8000 iteration, current mean episode reward is 0.385.
[INFO] In 9000 iteration, current mean episode reward is 0.161.
[INFO] In 10000 iteration, current mean episode reward is 0.325.
[INFO] In 11000 iteration, current mean episode reward is 0.423.
[INFO] In 12000 iteration, current mean episode reward is 0.539.
[INFO] In 13000 iteration, current mean episode reward is 0.216.
[INFO] In 14000 iteration, current mean episode reward is 0.476.
[INFO] In 15000 iteration, current mean episode reward is 0.547.
[INFO] In 16000 iteration, current mean episode reward is 0.53.
[INFO] In 17000 iteration, current mean episode reward is 0.449.
[INFO] In 18000 iteration, current mean episode reward is 0.3.
[INFO] In 19000 iteration, current mean episode reward is 0.647.
[INFO] In 20000 iteration, current mean episode reward is 0.434.
We expect to get the mean episode reward greater than 0.6. But you get: 0.4
34. Please check your codes.
```

```
In [47]: new_sarsa_trainer.render()
```

```
(Right)
SFFFFFFF
FFFFFFF
FFFHFFFF
FFFFFHFF
FFFHFFFF
FHHFFHF
FHFFHFHF
FFFHFFFG
```

Now you have implement the MC Control algorithm. You have finished this section. If you want to do more investigation like comparing the policy provided by SARSA, Q-Learning and MC Control, then you can do it in the next cells. It's OK to leave it blank.

```
In [48]: # Comparisons:

'''

1. SARSA vs Q-Learning

The only difference is the "next action choosing policy":
SARSA is epsilon-greedy while Q-Learning is absolute greedy.

Therefore, when convergence, SARSA's policy is more conservative,
while Q-Learning's policy is more bold and radical, and of course,
usually "better" (with less time/cost/distance).

'''

'''

2. Monte-Carlo vs TD

Episode-by-episode instead of step-by-step.

Theoretically, if there're enough episodes, MC will get the perfect solution,
hence the speed of MC is apparently slower than TD's.

'''

Out[48]: "\n\n2. Monte-Carlo vs TD\n\nEpisode-by-episode instead of step-by-step.\n\nTheoretically, if there're enough episodes, MC will get the perfect solution,\nhence the speed of MC is apparently slower than TD's.\n\n"
```

Conclusion and Discussion

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!

.

In []: