

▼ Mount the Google Drive onto the Colab as the storage location.

Following the instructions returned from the below cell. You will click a web link and select the google account you want to mount, then copy the authorization code to the blank, press enter.

```
1 # This must be run within a Google Colab environment
2 from google.colab import drive
3 drive.mount('/content/gdrive')
```

Mounted at /content/gdrive

▼ Append the directory location where you upload the start code folder (In this problem, *RLalgs*) to the sys.path

E.g. dir = '[/content/drive/My Drive/RL](#)/.', start code folder is inside "RL" folder.

```
1 import sys
2 sys.path.append('/content/gdrive/MyDrive/2021fall/e6885/coding_assignment/.')
3 # sys.path.append('</dir/to/start/code/folder/.>')
```

Your code should remain in the block marked by

```
#####
# YOUR CODE STARTS HERE
# YOUR CODE ENDS HERE
#####
```

Please don't edit anything outside the block.

```
1 %load_ext autoreload
2 %autoreload 2
```

The autoreload extension is already loaded. To reload it, use:
%reload_ext autoreload

```
1 import numpy as np
2 import random
3 import matplotlib.pyplot as plt
4 import gym
```

▼ 1. Incremental Implementation of Average

We've finished the incremental implementation of average for you. Please call the function estimate with 1/step step size and fixed step size to compare the difference between this two

on a simulated Bandit problem.

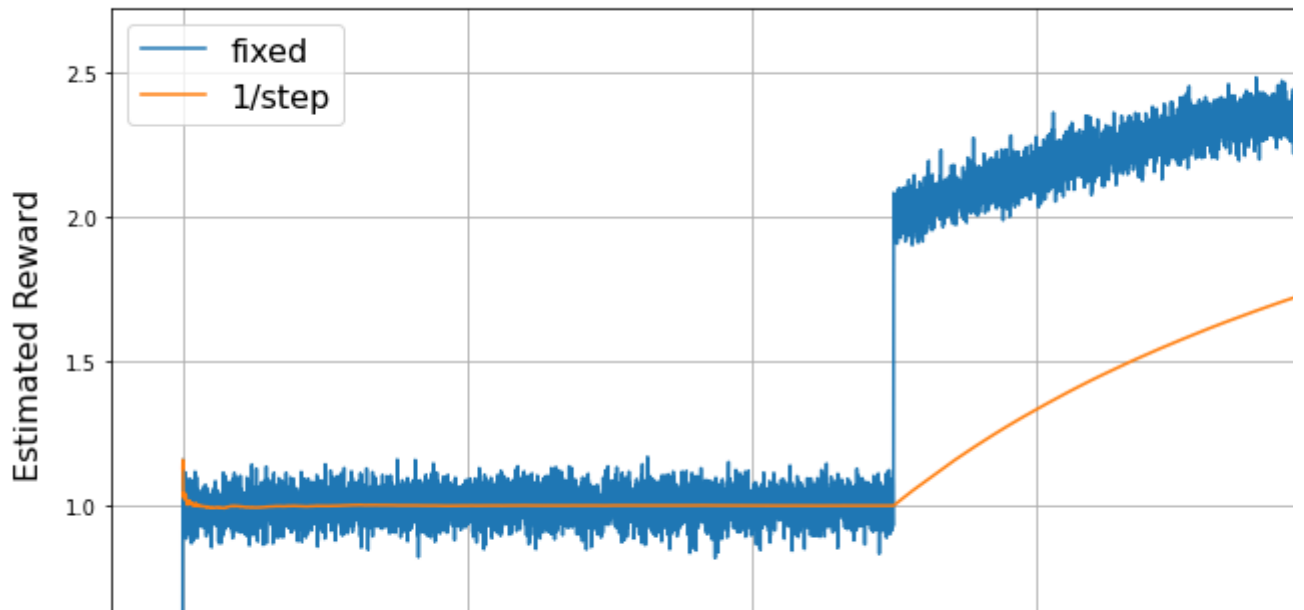
```
1 from RLalgs.utils import estimate
2 random.seed(6885)
3 numTimeStep = 10000
4 q_h = np.zeros(numTimeStep + 1) # Q Value estimate with 1/step step size
5 q_f = np.zeros(numTimeStep + 1) # Q value estimate with fixed step size
6 FixedStepSize = 0.5 #A large number to exaggerate the difference
7 for step in range(1, numTimeStep + 1):
8     if step < numTimeStep / 2:
9         # gauss distribution
10        r = random.gauss(mu = 1, sigma = 0.1)
11    else:
12        r = random.gauss(mu = 3, sigma = 0.1)
13
14    #TIPS: Call function estimate defined in ./RLalgs/utils.py
15    #####
16    # YOUR CODE STARTS HERE
17    q_h[step] = estimate(q_h[step-1], 1/step, r)
18    q_f[step] = estimate(q_h[step-1], FixedStepSize, r)
19    # YOUR CODE ENDS HERE
20    #####
21
22 q_h = q_h[1:]
23 q_f = q_f[1:]
```

RLalgs is a package containing Reinforcement Learning algorithms Epsilon-Greedy, Policy Iter:

Plot the two Q value estimates. (Please include a title, labels on both axes, and legends)

```
1 #####
2 # YOUR CODE STARTS HERE
3
4 plt.figure(figsize=(14,6))
5
6 # fixed stpsize
7 plt.plot(range(numTimeStep), q_f, label='fixed')
8 # 1/step stepsize
9 plt.plot(range(numTimeStep), q_h, label='1/step')
10 plt.grid()
11 plt.legend(fontsize=16)
12 plt.title('Estimated Reward of Bandit Problem', fontsize=24)
13 plt.ylabel('Estimated Reward', fontsize=16, labelpad=10)
14 plt.xlabel('Steps', fontsize=16, labelpad=10);
15
16 # YOUR CODE ENDS HERE
17 #####
```

Estimated Reward of Bandit Problem



▼ 2. ϵ -Greedy for Exploration

In Reinforcement Learning, we are always faced with the dilemma of exploration and exploitation. ϵ -Greedy is a trade-off between them. You are gonna implement Greedy and ϵ -Greedy. We combine these two policies in one function by treating Greedy as ϵ -Greedy where $\epsilon = 0$. Edit the function `epsilon_greedy` in [./RLalgs/utils.py](#).

```
1 from RLalgs.utils import epsilon_greedy
2 np.random.seed(6885) #Set the seed to cancel the randomness
3 q = np.random.normal(0, 1, size = 5)
4 #####
5 # YOUR CODE STARTS HERE
6 greedy_action = epsilon_greedy(q, 0) #Use epsilon = 0 for Greedy
7 e_greedy_action = epsilon_greedy(q, 0.1) #Use epsilon = 0.1
8 # YOUR CODE ENDS HERE
9 #####
10 print('Values:')
11 print(q)
12 print('Greedy Choice =', greedy_action)
13 print('Epsilon-Greedy Choice =', e_greedy_action)
```

```
Values:
[ 0.61264537  0.27923079 -0.84600857  0.05469574 -1.09990968]
Greedy Choice = 0
Epsilon-Greedy Choice = 0
```

You should get the following results.

Values:

```
[ 0.61264537  0.27923079 -0.84600857  0.05469574 -1.09990968]
```

Greedy Choice = 0

▼ 3. Frozen Lake Environment

```
1 env = gym.make('FrozenLake-v0')
```

▼ 3.1 Derive Q value from V value

Edit function `action_evaluation` in [./RLalgs/utils.py](#).

TIPS: $q(s, a) = \sum_{s', r} p(s', r | s, a)(r + \gamma v(s'))$

```
1 from RLalgs.utils import action_evaluation
2 v = np.ones(16)
3 q = action_evaluation(env = env.env, gamma = 1, v = v)
4 print('Action values:')
5 print(q)
```

```
Action values:
[[1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.      1.      1.      ]
 [1.      1.33333333 1.33333333 1.33333333]
 [1.      1.      1.      1.      ]]
```

You should get Q values all equal to one except at State 14

Pseudo-codes of the following four algorithms can be found on Page 80, 83, 130, 131 of the Sutton's book.

▼ 3.2 Model-based RL algorithms

```
1 from RLalgs.utils import action_evaluation, action_selection, render
```

▼ 3.2.1 Policy Iteration

Edit the function `policy_iteration` and relevant functions in [./RLalgs/pi.py](#) to implement the

```
1 from RLalgs.pi import policy_iteration
2 V, policy, numIterations = policy_iteration(env = env.env, gamma = 1, max_iteration =
3 print('State values:')
4 print(V)
5 print('Number of iterations to converge =', numIterations)

State values:
[0.82352769 0.82352712 0.82352673 0.82352653 0.82352786 0.
 0.52941059 0.          0.82352813 0.82352848 0.76470507 0.
 0.          0.88235229 0.94117614 0.          ]
Number of iterations to converge = 500
```

You should get values close to:

State values:

```
[0.82352774 0.8235272 0.82352682 0.82352662 0.82352791 0.
0.52941063 0. 0.82352817 0.82352851 0.76470509 0.
0. 0.88235232 0.94117615 0.]
```

```
1 #Uncomment and run the following to evaluate your result, comment them when you get
2 # Q = action_evaluation(env = env.env, gamma = 1, v = V)
3 # policy_estimate = action_selection(Q)
4 # render(env, policy_estimate)
```

▼ 3.2.2 Value Iteration

Edit the function `value_iteration` and relevant functions in [./RLalgs/vi.py](#) to implement the Value Iteration Algorithm.

```
1 from RLalgs.vi import value_iteration
2 V, policy, numIterations = value_iteration(env = env.env, gamma = 1, max_iteration =
3 print('State values:')
4 print(V)
5 print('Number of iterations to converge =', numIterations)

State values:
[0.82352513 0.82352369 0.82352267 0.82352214 0.82352544 0.
 0.5294087 0.          0.82352604 0.82352689 0.76470365 0.
 0.          0.88235115 0.94117554 0.          ]
Number of iterations to converge = 500
```

You should get values close to:

State values:

```
[0.82352773 0.82352718 0.8235268 0.8235266 0.8235279 0.
0.52941062 0. 0.82352816 0.8235285 0.76470509 0.
0. 0.88235231 0.94117615 0.]
```

```

1 #Uncomment and run the following to evaluate your result, comment them when you get
2 # Q = action_evaluation(env = env.env, gamma = 1, v = V)
3 # policy_estimate = action_selection(Q)
4 # render(env, policy_estimate)

```

3.3 Model free RL algorithms

▼ 3.3.1 Q-Learning

Edit the function QLearning in [./RLalgs/ql.py](#) to implement the Q-Learning Algorithm.

```

1 from RLalgs.ql import QLearning
2 Q = QLearning(env = env.env, num_episodes = 1000, gamma = 1, lr = 0.1, e = 0.1)
3 print('Action values:')
4 print(Q)

```

```

Action values:
[[2.11199210e-03 3.67909924e-02 3.60097180e-03 1.08141550e-02]
 [4.47010439e-03 2.63382467e-02 3.84696543e-02 1.33505818e-02]
 [8.13408307e-02 2.51285268e-02 3.30812594e-02 1.18327649e-02]
 [2.51288321e-02 3.10232495e-03 7.19206088e-05 4.24684003e-05]
 [6.55202187e-02 4.36773528e-03 8.78594226e-03 5.15825836e-03]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [8.56434774e-02 2.97735169e-02 4.43788385e-02 2.83638073e-03]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [8.00560910e-03 9.28663313e-02 3.80966177e-02 3.16427834e-02]
 [0.00000000e+00 2.39119991e-01 1.11204636e-01 7.51017004e-02]
 [1.09846695e-01 8.27005453e-02 2.97020742e-01 8.02628561e-03]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [6.73159352e-02 9.08349828e-02 7.21099186e-02 2.14149757e-01]
 [1.78006566e-01 6.54237672e-01 3.21100169e-01 2.51990747e-01]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]]

```

Generally, you should get non-zero action values on non-terminal states.

```

1 #Uncomment the following to evaluate your result, comment them when you generate the
2 # env = gym.make('FrozenLake-v1')
3 # policy_estimate = action_selection(Q)
4 # render(env, policy_estimate)

```

▼ 3.3.2 SARSA

Edit the function SARSA in [./RLalgs/sarsa.py](#) to implement the SARSA Algorithm.

```

1 from RLalgs.sarsa import SARSA
2 Q = SARSA(env = env.env, num_episodes = 1000, gamma = 1, lr = 0.1, e = 0.1)
3 print('Action values:')
4 print(Q)

```

Action values:

```
[[7.46111235e-02 6.73376815e-02 1.08750153e-01 5.56249585e-02]
 [3.07835188e-02 2.18086590e-02 2.41144029e-02 1.13791892e-01]
 [1.48316881e-01 2.25421919e-02 6.39042944e-02 3.25353974e-02]
 [2.07363577e-02 9.38229063e-03 1.03752717e-02 7.68844233e-02]
 [1.38554506e-01 9.82300893e-02 3.70805732e-02 5.50829298e-02]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [1.48074228e-01 2.01559919e-02 5.49124151e-02 1.81843732e-04]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [8.18453622e-02 1.21912079e-01 5.07193900e-02 1.81406611e-01]
 [5.73904157e-02 2.33346024e-01 1.31252032e-01 4.75127127e-02]
 [3.14899445e-01 2.88634098e-01 2.13909429e-01 5.06246406e-02]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]
 [1.17805420e-01 2.08277949e-01 1.68114919e-01 1.52920155e-01]
 [2.21025847e-01 5.08879901e-01 3.64096881e-01 3.97002465e-01]
 [0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00]]
```

Generally, you should get non-zero action values on non-terminal states.

```
1 #Uncomment the following to evaluate your result, comment them when you generate the
2 # env = gym.make('FrozenLake-v0')
3 # policy_estimate = action_selection(Q)
4 # render(env, policy_estimate)
```

▼ 3.4 Human

You can play this game if you are interested. See if you can get the frisbee either with or without the model.

```
1 from RLalgs.utils import human_play
2 # Uncomment and run the following to play the game, comment it when you generate
3 # env = gym.make('FrozenLake-v1')
4 # human_play(env)
```

▼ 4. Exploration VS. Exploitation

Try to reproduce Figure 2.2 (the upper one is enough) of the Sutton's book based on the experiment described in [Chapter 2.3](#).

```
1 from RLalgs.utils import *
2
3 class Bandit():
4     def __init__(self, num_arm, num_iter, eps):
5         """
6         Bandit problem following e-greedy policy, the arm machine follow normal
7
8         Inputs:
9         - num_arm: the number of arm machine
```

```

10         - num_iter: the number of total actions to be taken
11         - eps: epsilon in e-greedy policy
12
13         """
14         # number of arms
15         self.num_arm = num_arm
16         # number of total actions to be taken
17         self.num_iter = num_iter
18         # epsilon
19         self.eps = eps
20         # record average reward for the whole game
21         self.reward = 0
22         # self.reward_hist = [0]
23         self.reward_hist = np.zeros(num_iter+1)
24         self.reward_hist[0] = 0
25         # record average reward for each arms
26         self.q = np.zeros(num_arm)
27         # number of actions taken for the whole game
28         self.num_step = 0
29         # number of actions taken on each arms
30         self.num_step_arm = np.zeros(num_arm)
31         # initialize arms
32         self.arms = None
33         self.init_arms()
34
35     def reset(self, eps, num_iter):
36         """
37         reset eps and reward
38         """
39         self.reward = 0
40         # self.reward_hist = [0]
41         self.reward_hist = np.zeros(num_iter+1)
42         self.reward_hist[0] = 0
43         self.q = np.zeros_like(self.q)
44         self.num_step = 0
45         self.num_step_arm = np.zeros_like(self.num_step_arm)
46         self.eps = eps
47         self.num_iter = num_iter
48
49     def init_arms(self):
50         """
51         Initialize arms by randomly choose mu and std for normal distribution
52
53         mu follows normal distribution
54         std = 1
55
56         store in array self.arms, which has size 2*num_arm
57         first row of arms is mu
58         second row of arms is std
59         """
60         assert self.arms == None, 'Arms already defined'
61
62         # first row is mu, secound row is std
63         # arms = np.random.rand(2, self.num_arm)
64         arms = np.random.normal(0, 1, (2, self.num_arm))

```



```

65
66         # arms[0,:] = arms[0,:]*4 - 2
67         # arms[1,:] = arms[1,:]*10 + 20
68         arms[1,:] = 1
69
70         self.arms = arms
71
72     def take_action(self, a):
73         """
74         Take action a, range of a is [0,num_arm-1], return the reward of the
75
76         Input:
77         - a: the action to be taken
78
79         Output:
80         - r: reward of that action
81         """
82
83         return np.random.normal(self.arms[0,a], self.arms[1,a])
84
85     def step(self):
86         """
87         one step forward, choose the action following the e-greedy policy
88         """
89
90         # e-greedy policy
91         a = epsilon_greedy(self.q, self.eps)
92         # get reward of that action
93         r = self.take_action(a)
94         # update steps taken
95         self.num_step += 1
96         self.num_step_arm[a] += 1
97         # update the total average reward and the average reward of that arm
98         self.reward = estimate(self.reward, 1/self.num_step, r)
99         # self.reward_hist.append(self.reward)
100        self.q[a] = estimate(self.q[a], 1/self.num_step_arm[a], r)
101
102    def run(self):
103        """
104        play the game for iteration times
105
106        Output:
107        - history of total average reward
108        """
109        for i in range(self.num_iter):
110            self.step()
111            self.reward_hist[i+1] = self.reward
112
113        return self.reward_hist
114
115
116

```

You should get curves similar to that in the book.

```

1 # Plot the average reward
2 #####
3
4
5 num_arms = 10
6 num_iter = 1000
7 num_repeat = 1000
8
9 hist_0 = np.zeros(num_iter+1)
10 hist_001 = np.zeros(num_iter+1)
11 hist_01 = np.zeros(num_iter+1)
12
13
14 for i in range(1, num_repeat+1):
15     eps = 0
16     bandit = Bandit(num_arms, num_iter, 0)
17     hist_0 += bandit.run()
18
19     # eps = 0.01
20     bandit.reset(0.01, num_iter)
21     hist_001 += bandit.run()
22
23
24     # eps = 0.1
25     bandit.reset(0.1, num_iter)
26     hist_01 += bandit.run()
27
28     if i%100 == 0:
29         print('Repeated %d times' % (i))
30
31 hist_0 /= num_repeat
32 hist_001 /= num_repeat
33 hist_01 /= num_repeat
34
35 plt.figure()
36 plt.plot(hist_0, 'g', label="$\epsilon=0$ (Greedy)")
37 plt.plot(hist_001, 'r', label="$\epsilon=0.01$")
38 plt.plot(hist_01, 'b', label="$\epsilon=0.1$")
39 plt.title("Average $\epsilon$-greedy rewards after %d Repeats" % (num_repeat))
40 plt.xlabel("Steps")
41 plt.ylabel("Average Reward")
42 plt.legend();
43
44
45 # YOUR CODE ENDS HERE
46 #####

```

Repeated 100 times
Repeated 200 times
Repeated 300 times
Repeated 400 times
Repeated 500 times
Repeated 600 times
Repeated 700 times
Repeated 800 times
Repeated 900 times
Repeated 1000 times

