# Exam 2: Advanced AI

Caleb Ehrisman

11/30/2022

## Overview

The objective of this assignment is to implement a Fourier Basis SARSA($\lambda$) algorithm to solve the mountain car problem.

The mountain car is a popular problem within Reinforcement Learning where an under-powered car is stuck is a valley. To get out the car must build momentum by accelerating left and right until it gains enough to reach the peak.
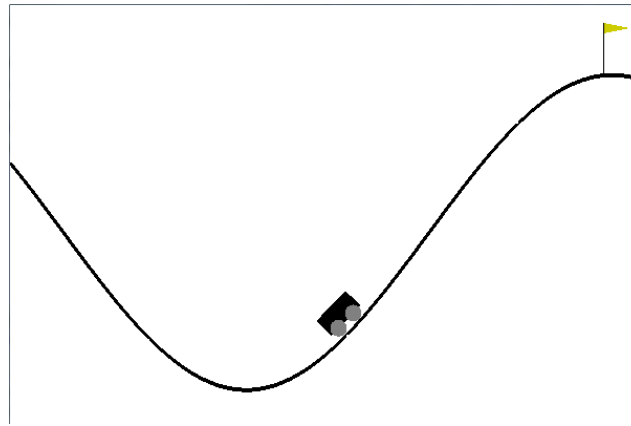


Figure 1: Screenshot of a graphic of mountain car

## Approach and Implementation

The first step in the approach was to make sure I always had a beer in hand to ensure proper usage of Balmer's Peak. In the solution implementation, True Online SARSA($\lambda$) was used as the learning agent.



**True online Sarsa($\lambda$) for estimating $\mathbf{w}^\top \mathbf{x} \approx q_\pi$ or $q_*$**

Input: a feature function $\mathbf{x} : \mathcal{S}^+ \times \mathcal{A} \to \mathbb{R}^d$ such that $\mathbf{x}(terminal, \cdot) = \mathbf{0}$
Input: a policy $\pi$ (if estimating $q_\pi$)
Algorithm parameters: step size $\alpha > 0$, trace decay rate $\lambda \in [0, 1]$, small $\varepsilon > 0$
Initialize: $\mathbf{w} \in \mathbb{R}^d$ (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:
    Initialize $S$
    Choose $A \sim \pi(\cdot|S)$ or $\varepsilon$-greedy according to $\hat{q}(S, \cdot, \mathbf{w})$
    $\mathbf{x} \leftarrow \mathbf{x}(S, A)$
    $\mathbf{z} \leftarrow \mathbf{0}$
    $Q_{old} \leftarrow 0$
    Loop for each step of episode:
    |    Take action $A$, observe $R$, $S'$
    |    Choose $A' \sim \pi(\cdot|S')$ or $\varepsilon$-greedy according to $\hat{q}(S', \cdot, \mathbf{w})$
    |    $\mathbf{x}' \leftarrow \mathbf{x}(S', A')$
    |    $Q \leftarrow \mathbf{w}^\top \mathbf{x}$
    |    $Q' \leftarrow \mathbf{w}^\top \mathbf{x}'$
    |    $\delta \leftarrow R + \gamma Q' - Q$
    |    $\mathbf{z} \leftarrow \gamma\lambda\mathbf{z} + \left(1 - \alpha\gamma\lambda\mathbf{z}^\top\mathbf{x}\right)\mathbf{x}$
    |    $\mathbf{w} \leftarrow \mathbf{w} + \alpha(\delta + Q - Q_{old})\mathbf{z} - \alpha(Q - Q_{old})\mathbf{x}$
    |    $Q_{old} \leftarrow Q'$
    |    $\mathbf{x} \leftarrow \mathbf{x}'$
    |    $A \leftarrow A'$
    until $S'$ is terminal

Figure 2: SARSA($\lambda$) from Sutton and Barto RL book

Where $\mathbf{w}$ is the weight vector and $\mathbf{x}$ is the feature approximation vector for a given state and action. The $\mathbf{z}$ is the trace vector. The hyper parameters are as follows: $\alpha$ - learning rate, $\lambda$ - decay rate of trace, and $\epsilon$ - value to choose actions based on probability.

For the approximator, the one from the paper linked in the assignment ,Value Function Approximation in Reinforcement Learning using the Fourier Basis, was implemented to estimate the Q values for a state-action pairs. Using N(dimensions of the state space) and the M(order of the Fourier function) to create a matrix of $(M+1)^2$ x $N$ for the co-efficients. To get the values of this matrix it is required to call get_features within the

FourierBasis class. This function uses equation (1) as shown in the book.

$$\phi_i(x) = cos(\pi \mathbf{c}^i \cdot \mathbf{x}) \tag{1}$$

## Results

The program has several options to run from the terminal all handled through argparse. The optional arguments are:

- --render
  type=str   True if want to run with graphics.

- --order
  type=int   default=3, Choose any order for the fourier basis

- --num_epochs
  type=int   default=1000, Choose any number of epochs to run

- --fourier
  type=str   default=True, if false then run standard SARSA($\lambda$) without Fourier basis

- --file
  type=str   default='weights.npy', File path to save the weights after training. Must be a '.npy' file.

- --train
  type=str   default=True, If false will load file and run without updates

- --eval
  type=str   default=False, Will run multiple bases and plot

Example usage:     C:\> py  main.py  --render  True  --order  5  --num_epoch  5000

With the file loading and saving it makes it really easy to train the weights and save the values into a file. The run method will load these values and then allow you to simply run a greedy model on the weights. If the weights were trained properly then it should quickly reach the terminal state and consistently do so.
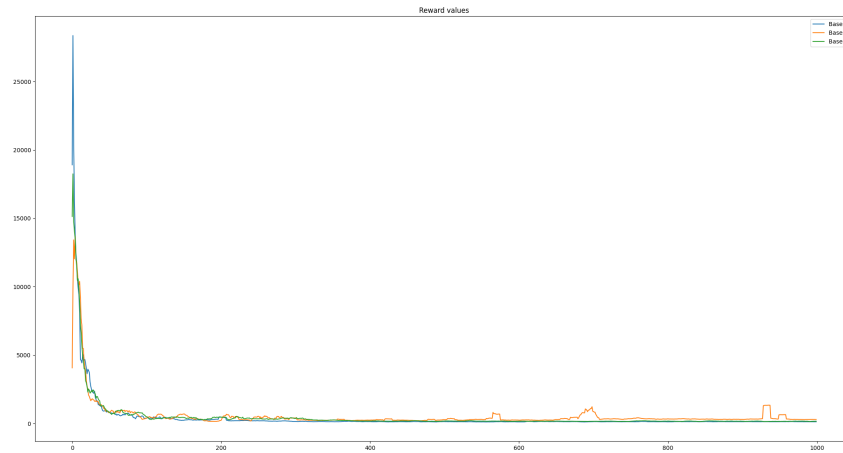
Figure 3: Learning Curves of bases 3,5,7
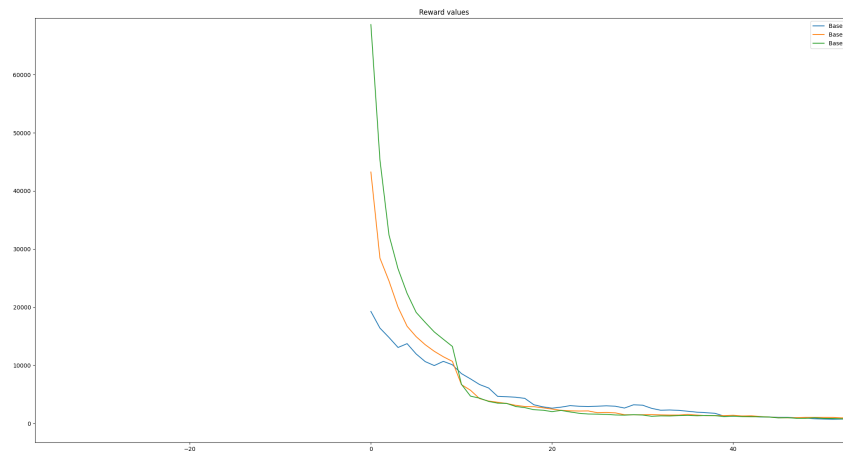


Figure 4: The first 50 episodes

Figs 3, 4 are the learning curves of the different Fourier bases. The lines are batched averages of 10 episodes each to present the lines in a smooth way while still clearly seeing how the curve is trending. It appears that it converges roughly 50-100 episodes into training with values of $\alpha$ - 0.01, $\lambda$ - 0.9, and $\epsilon$ - 0.05.

The surface plots of the Q_value estimation form the Fourier Basis Approximation
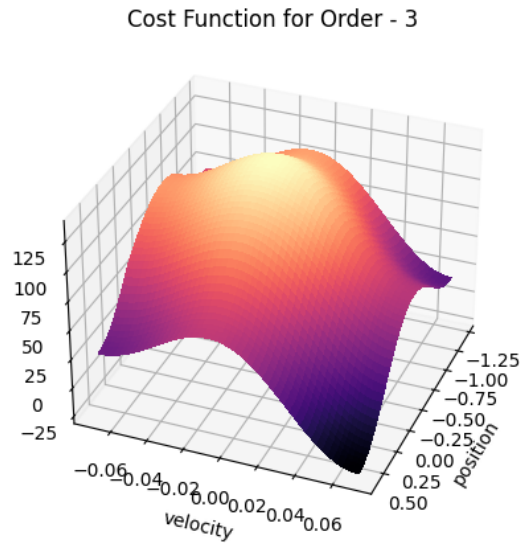


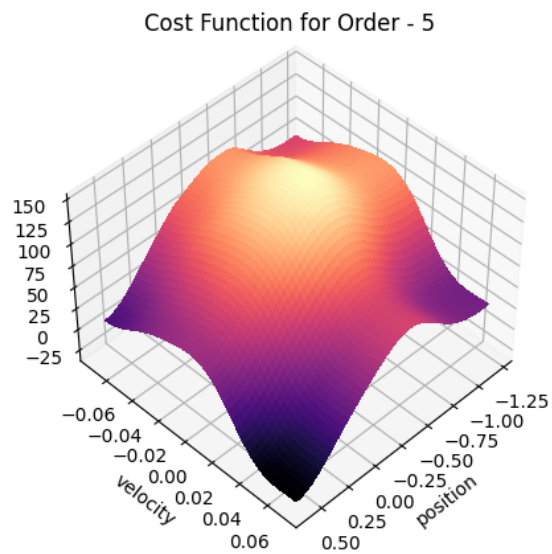Figure 5: Surface plot of base 3



Figure 6: Surface plot of base 5

Figure 7: Surface plot of base 7

## Question

**The Mountain Car contains a negative step reward and a zero goal reward. What would happen if $\gamma$ was less than 1 and the solution was many steps long?**

To test what happens when the gamma is less than 1.0, simply just change the gamma value and let it run for infinite steps. As can be seen in Fig 8. the learning curve is interesting. The gamma seems to pose an issue where some runs will end up being absolutely horrible. It will still converge; however, the solution is still unstable as can be seen in Fig 9 with some runs going to 40,000 or more. With the gamma being 1.0 the worst runs still hardly reach over 20,000 steps, so reducing the gamma below 1.0 has negative effects on the training.

Figure 8: Learning Curves basis 3 and gamma 0.95



Figure 9: Last 200 episodes.

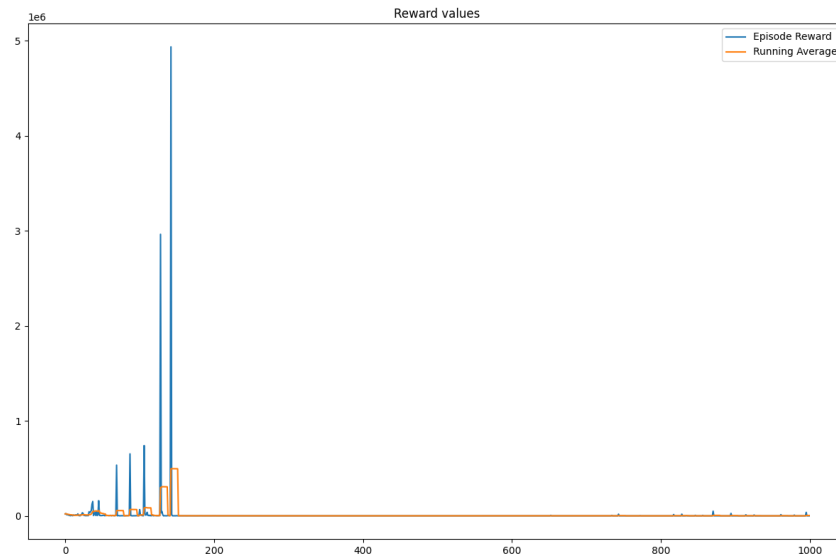**What would happen if we had a zero step cost and a positive goal reward, for the case where** $\gamma$ **= 1, and the case where** $\gamma$ **< 1?** In the case of just changing the time step reward from -1 to 0 and the terminal state reward to 1 instead of 0 increases the training time significantly and it struggles to converge (took over an hour to reach 600 epochs). When setting the reward to 100 for the terminal state it was able to learn a solution early on but then would diverge later on as can be seen in Fig 10 and Fig 11. And the early solution was not close to optimal. The issue is that it is not getting punished for taking bad routes and sometimes it creates a path that just loops on itself and it will not correct due to no bad reward. So trying to test thing through code is a hassle and very unstable. Some runs will actually converge and others will not and seem to fall in a loop as it will sit on an episode for upwards of 30 min. And changing the gamma value will not remedy this as we saw in the Fig. 10 it will generally get worse with more episodes.



Figure 10: Learning Curves basis 3 and 0 step reward and 100 terminal state reward

Figure 11: Early learning curve

## Final Remarks

The github repo that hosts the code is at GitHub

The code is also included in a zip submitted with the pdf and below in the appendix. There is a default weight vector stored in 'default.npy'. If desired to run then do　　C:\> py main.py --render True --file default

## 1　Appendix

```
1  '''
2  Author: Caleb Ehrisman
3  Course- Advanced AI CSC-549
4  Assignment - Programming Assignment #3 - Mountain Car
5
6  Tasks
7   - Implement Sarsa(lambda) to solve mountain car problem
8   - Use Linear Function Approximation with Fourier Basis
        functions
9   - Show different learning curves for 3rd, 5th, and 7th
        order Fourier bases
10  - Create surface plot of the value function
```

```python
11   - Answer short response question
12  '''
13
14  import gym
15  import matplotlib
16
17  from agent import Agent
18  import argparse
19  import pandas as pd
20  import matplotlib.pyplot as plt
21  import numpy as np
22
23  """
24      parse gathers command line arguments.
25
26
27      :return: a list of all parsed arguments
28      """
29
30
31  def parse():
32      parser = argparse.ArgumentParser()
33      parser.add_argument('--render', type=str, help='Specify
            to run simulation or not')
34      parser.add_argument('--order', type=int, help='Choose
            order for fourier basis', default=3)
35      parser.add_argument('--num_epochs', type=int, help='
            Choose number of epochs', default=1000)
36      parser.add_argument('--fourier', type=str, help='Choose
            to use fourier', default=True)
37      parser.add_argument('--file', type=str,
38                          help='File path to save weights to.
                                Must be given with .npy extension
                                ', default='weights.npy')
39      parser.add_argument('--train', type=str, help='Choose if
             training or running', default='True')
40      parser.add_argument('--eval', type=str, default='False')
41      return parser.parse_args()
42
43
44  if __name__ == "__main__":
45      args = parse()
```

```python
46
47     if args.render == "True":
48         env = gym.make("MountainCar-v0", render_mode="human"
               )
49
50     else:
51         env = gym.make("MountainCar-v0")
52
53     file = args.file
54
55     n = args.num_epochs
56     if args.fourier == 'False':
57         agent = Agent(env, file, fourier=False, order=3)
58     else:
59         agent = Agent(env, file, order=3)
60
61     if args.eval == 'False':
62         if args.train == 'True':
63             rewards, avg, learner = agent.learn(n)
64         else:
65             rewards, avg = agent.run(n)
66
67         fig, ax = plt.subplots(figsize=(10, 4))
68         plt.plot(np.negative(rewards), label='Episode Reward
               ')
69         plt.plot(np.negative(avg), label='Running Average')
70         ax.set_title("Reward values")
71         plt.legend()
72         plt.show()
73
74         rewards = []
75         base = [3, 5, ]
76
77         rewards.append(avg)
78
79         fig, ax = plt.subplots(figsize=(10, 4))
80         plt.plot(np.negative(rewards[0]), label='Base 3')
81         plt.plot(np.negative(rewards[1]), label='Base 5')
82         plt.plot(np.negative(rewards[2]), label='Base 7')
83         ax.set_title("Reward values")
84         plt.legend()
85         plt.show()
```

```
86
87            low = env.observation_space.low
88            high = env.observation_space.high
89            difference = high - low
90
91            x_axis = np.linspace(low[0], high[0])
92            y_axis = np.linspace(low[1], high[1])
93            x_axis, y_axis = np.meshgrid(x_axis, y_axis)
94            z_axis = np.zeros(x_axis.shape)
95
96            for i in range(0, z_axis.shape[0]):
97                for j in range(0, z_axis.shape[1]):
98                    s = [(x_axis[i, j] - low[0]) / (high[0] -
                            low[0]), (y_axis[i, j] - low[1]) / (high
                            [1] - low[1])]
99                    (zq, _) = learner.best_action(s)
100                   z_axis[i, j] = -1.0 * zq
101
102               fig = plt.figure()
103               ax = plt.axes(projection='3d')
104               ax.plot_surface(x_axis, y_axis, z_axis, cmap=
                        matplotlib.cm.get_cmap("magma"))
105               ax.set_xlabel('position')
106               ax.set_ylabel('velocity')
107               ax.set_title('Cost Function for Order - ' + str(
                        n))
108               plt.show()
109
110    else:
111        rewards = []
112        base = [3, 5, 7]
113        learner = []
114        for i in range(3):
115            agent = Agent(env, file, order=base[i])
116            reward, avg, temp_learner = agent.learn(1000)
117            rewards.append(avg)
118            learner.append(temp_learner)
119
120        fig, ax = plt.subplots(figsize=(10, 4))
121        plt.plot(np.negative(rewards[0]), label='Base 3')
122        plt.plot(np.negative(rewards[1]), label='Base 5')
123        plt.plot(np.negative(rewards[2]), label='Base 7')
```

```
124            ax.set_title("Reward values")
125            plt.legend()
126            plt.show()
127
128            low = env.observation_space.low
129            high = env.observation_space.high
130            difference = high - low
131
132            x_axis = np.linspace(low[0], high[0])
133            y_axis = np.linspace(low[1], high[1])
134            x_axis, y_axis = np.meshgrid(x_axis, y_axis)
135            z_axis = np.zeros(x_axis.shape)
136
137            for b in range(3):
138                for i in range(0, z_axis.shape[0]):
139                    for j in range(0, z_axis.shape[1]):
140                        s = [(x_axis[i, j] - low[0]) / (high[0]
                                - low[0]), (y_axis[i, j] - low[1]) /
                                (high[1] - low[1])]
141                        (zq, _) = learner[b].best_action(s)
142                        z_axis[i, j] = -1.0 * zq
143
144                fig = plt.figure()
145                ax = plt.axes(projection='3d')
146                ax.plot_surface(x_axis, y_axis, z_axis, cmap=
                       matplotlib.cm.get_cmap("magma"))
147                ax.set_xlabel('position')
148                ax.set_ylabel('velocity')
149                ax.set_title('Cost Function for Order - ' + str(
                       base[b]))
150                plt.show()
151                '''
152 Author: Caleb Ehrisman
153 Course- Advanced AI CSC-549
154 Assignment - Programming Assignment #3 - Mountain Car
155
156 This file contains the code to implement the SARSA(lambda)
       algorithm.
157
158 All functions needed by solely the agent are included as
       member functions of class Agent
159 '''
```

```
160  import numpy as np
161  from fourier_basis import FourierBasis
162  from sarsalambdaFA import SarsaLambdaFA
163  from sarsa import Sarsa
164  import os.path
165
166  ALPHA = 0.0001
167  GAMMA = 1
168  EPSILON = 0.5
169  LAMBDA = 0.9
170
171
172  class Agent:
173
174      def __init__(self, environment, file, fourier=True,
             order=3, runs=1, gamma=0.001):
175          """
176                  init is the constructor for the Agent class.
177
178                  :param environment
179                  :return None
180          """
181          self.runs = runs
182          self.order = order
183          self.env = environment
184          self.gamma = gamma
185          self.num_actions = self.env.action_space.n
186          self.state_dims = self.env.observation_space.shape
                [0]
187          self.fourier = fourier
188          self.epoch_rewards = []
189          self.epoch_rewards_table = {'ep': [], 'avg': [], '
                min': [], 'max': []}
190          self.epoch_max_pos = []
191          self.file = file
192
193      def learn(self, num_epochs):
194          """
195                  Agent.learn does the actual stepping through and
                    exploring the environment and then updates
                    the Q_table if
```

```
196                 using traditional SARSA and updates the weight
                        and lambda vectors is using a fourier basis
197
198                 :param num_epochs
199                 :return None
200                 """
201         for run in range(0, self.runs):
202             fb = FourierBasis(state_space=self.env.
                    observation_space.shape[0], order=self.order)
203             if self.fourier:
204                 learner = SarsaLambdaFA(fa=fb, num_actions=
                        self.num_actions, alpha=0.0001, epsilon
                        =0.8)
205             else:
206                 learner = Sarsa(environment=self.env)
207
208             for i in range(num_epochs):
209
210                 learner.epsilon *= .99
211                 if self.fourier:
212                     learner.z = np.zeros(learner.w.shape)
213                 state, _ = self.env.reset()
214                 if self.fourier:
215                     state = (state - self.env.
                            observation_space.low) / (self.env.
                            observation_space.high - self.env.
                            observation_space.low)
216                 else:
217                     state = learner.discretized_env_state(
                            state)
218                     learner.E_table = learner.create_q_table
                            ()
219                 # steps = 0
220                 action = learner.action(state)
221                 done = False
222                 reward_sum = 0
223
224                 while not done:
225                     next_state, reward, done, info, _ = self
                            .env.step(action)
226                     # reward += 1
227                     # if done:
```

```
228                              #       reward = 100
229                          if self.fourier:
230                              next_state = (next_state - self.env.
                                     observation_space.low) / (
231                                        self.env.
                                             observation_space.
                                             high - self.env.
                                             observation_space.low
                                             )
232                          else:
233                              next_state = learner.
                                     discretized_env_state(next_state)
234
235                          next_action = learner.action(next_state)
236
237                          learner.update(state, action, reward,
                                next_state, done, next_action)
238
239                          # steps += 1
240                          state = next_state
241                          action = next_action
242                          reward_sum += reward
243
244                      #  Append max position data and reward data
                            for evaluation
245                      self.epoch_rewards.append(reward_sum)
246
247                      self.terminal_output(i)
248          np.save(self.file, learner.w)
249          return self.epoch_rewards, self.epoch_rewards_table[
                'avg'], learner
250
251      def run(self, num_epochs):
252          """
253              Agent.run uses a pre-trained set of weights to
                    greedily choose actions.
254
255              :param num_epochs
256              :return None
257          """
258
259          if os.path.exists(self.file):
```

```
260                 w = np.load(self.file)
261             else:
262                 print("Error loading file. Not found.")
263                 return
264
265         fb = FourierBasis(state_space=self.env.
                 observation_space.shape[0], order=self.order)
266         learner = SarsaLambdaFA(fa=fb, num_actions=self.
                 num_actions, alpha=0.0, epsilon=0.0)
267         learner.w = w
268
269         for i in range(num_epochs):
270             state, _ = self.env.reset()
271
272             state = (state - self.env.observation_space.low)
                     / (
273                         self.env.observation_space.high -
                             self.env.observation_space.low)
274
275             action = learner.action(state)
276             done = False
277             reward_sum = 0
278
279             while not done:
280                 next_state, reward, done, info, _ = self.env
                         .step(action)
281                 print(reward)
282                 next_state = (next_state - self.env.
                     observation_space.low) / (
283                         self.env.observation_space.high -
                             self.env.observation_space.low)
284
285                 next_action = learner.action(next_state)
286                 action = next_action
287                 reward_sum += reward
288
289             #  Append max position data and reward data for
                 evaluation
290             self.epoch_rewards.append(reward_sum)
291
292             self.terminal_output(i)
293
```

```python
294             return self.epoch_rewards, self.epoch_rewards_table[
                    'avg']
295
296     def terminal_output(self, i):
297         # Terminal Output for stats of each epoch
298         avg_reward = sum(self.epoch_rewards[-10:]) / len(
                self.epoch_rewards[-10:])
299         self.epoch_rewards_table['ep'].append(i)
300         self.epoch_rewards_table['avg'].append(avg_reward)
301         self.epoch_rewards_table['min'].append(min(self.
                epoch_rewards[:]))
302         self.epoch_rewards_table['max'].append(max(self.
                epoch_rewards[:]))
303
304         print(f"Epoch - {i}\t| avg: {avg_reward:.2f}\t| min:
                {min(self.epoch_rewards[-1:]):.2f}"
305               f"\t| max: {max(self.epoch_rewards[-1:]):.2f}"
                  )
306
307 '''
308 Author: Caleb Ehrisman
309 Course- Advanced AI CSC-549
310 Assignment - Programming Assignment #3 - Mountain Car
311
312 This file contains the code to implement the SARSA(lambda)
    with a function approximator.
313
314 '''
315 class SarsaLambdaFA:
316     def __init__(self, fa, num_actions=None, alpha=0.01,
            gamma=1.0, lamb=0.9, epsilon=0.5):
317         self.gamma = gamma
318         self.lamb = lamb
319         self.epsilon = epsilon
320         self.alpha = alpha
321         self.num_actions = num_actions
322         self.fourier_basis = []
323
324         for i in range(0, self.num_actions):
325             self.fourier_basis.append(copy.deepcopy(fa))
326
```

```python
327             self.w = np.zeros([self.fourier_basis[0].coeff.shape
                    [0], num_actions])
328             self.z = np.zeros(self.w.shape)
329
330             self.w[0, :] = 0.0
331
332         def action(self, state):
333             """
334                     Agent.action determines what action to take
                            based on state
335
336                     :param state:
337                     :return action
338             """
339             if np.random.uniform(0, 1) < self.epsilon:
340                 return random.randrange(0, self.num_actions)
341
342             best = -math.inf
343             best_actions = []
344             for a in range(0, self.num_actions):
345                 q = self.Q(state, a)
346                 if math.isclose(q, best):
347                     best_actions.append(a)
348                 elif q > best:
349                     best = q
350                     best_actions = [a]
351
352             return random.choice(best_actions)
353
354         def Q(self, state, action):
355             return np.dot(self.w[:, action], self.fourier_basis[
                    action].get_features(state))
356
357         def best_action(self, state):
358
359             best = -math.inf
360             best_action = 0
361             for a in range(0, self.num_actions):
362                 q = self.Q(state, a)
363                 if q > best:
364                     best = q
365                     best_action = a
```

```python
366            return best, best_action
367
368    def update(self, state, action, reward, next_state, done
           , next_action=None):
369        """
370            Agent.update updates the Q table based on the
                   SARSA algorithm. It also updates the trace
                   table
371
372            :param next_action:
373            :param done:
374            :param next_state:
375            :param state:
376            :param action:
377            :param reward
378            :return None
379        """
380
381        delta = reward - self.Q(state, action)
382
383        if not done:
384            if next_action is not None:
385                delta += self.gamma * self.Q(next_state,
                       next_action)
386            else:
387                q_dot, next_action = self.best_action(
                       next_state)
388                delta += self.gamma * self.best_action(q_dot
                       )
389
390        phi = self.fourier_basis[action].get_features(state)
391
392        for a in range(0, self.num_actions):
393            self.z[:, a] *= self.gamma * self.lamb
394            if a == action:
395                self.z[:, a] += phi
396            self.w[:, a] += self.alpha * delta * self.z[:, a
                   ]
397
398        return delta
399
400        '''
```

```
401  Author: Caleb Ehrisman
402  Course- Advanced AI CSC-549
403  Assignment - Programming Assignment #3 - Mountain Car
404  This file contains the code to implement the SARSA(lambda)
         algorithm.
405  All functions needed by solely the agent are included as
         member functions of class Agent
406  '''
407  import numpy as np
408
409  ALPHA = 0.1
410  GAMMA = 1
411  EPSILON = 0.5
412  LAMBDA = 0.9
413
414
415  class Sarsa:
416      def __init__(self, environment, gamma=1.0, epsilon=0.5,
             alpha=0.0001, lamb=0.9):
417          """
418                  init is the constructor for the Agent class.
419                  :param environment
420                  :return None
421          """
422          self.env = environment
423          self.E_table = self.create_q_table()
424          self.Q_table = self.create_q_table()
425          self.alpha = alpha
426          self.gamma = gamma
427          self.epsilon = epsilon
428          self.lamb = lamb
429          self.epoch_rewards = []
430          self.epoch_rewards_table = {'ep': [], 'avg': [], '
                 min': [], 'max': []}
431          self.epoch_max_pos = []
432
433      def create_q_table(self):
434          """
435                  Agent.create_q_table creates the Q table
                     that fits all states
436                  :return np.array of [x_lim][y_lim][
                     num_actions]
```

```
437                """
438                high = self.env.observation_space.high
439                low = self.env.observation_space.low
440                num_states = (high - low) * np.array([10, 100])
441                num_states = np.round(num_states, 0).astype(int) + 1
442                num_actions = self.env.action_space.n
443                return np.zeros([num_states[0], num_states[1],
                       num_actions])
444
445         def action(self, state):
446                """
447                        Agent.action determines what action to take
                           based on state
448                        :param state
449                        :return action
450                """
451                if np.random.uniform(0, 1) < EPSILON:
452                    action = self.env.action_space.sample()
453                else:
454                    # disc_state = self.discretized_env_state(state)
455                    action = np.argmax(self.Q_table[state[0], state
                          [1]])
456                return action
457
458         def update(self, state, action, reward, next_state, done
                , next_action):
459                """
460                    Agent.update updates the Q table based on the
                        SARSA algorithm. It also updates the trace
                        table
461                    :param done:
462                    :param state
463                    :param action
464                    :param reward
465                    :param next_state
466                    :param next_action
467                    :return None
468                """
469                target = reward + self.gamma * self.Q_table[
                       next_state[0], next_state[1], next_action]
470
```

```
471             error = target - self.Q_table[state[0], state[1],
                    action]
472             # print(self.E_table[state[0], state[1], action])
473             self.E_table[state[0], state[1], action] += 1
474
475             self.Q_table += 0.01 * error * self.E_table
476
477             self.E_table *= self.gamma * self.lamb
478
479     def discretized_env_state(self, state):
480         """
481             Agent.discretized_env_state takes a given state
                    and discretizes the state to use whole
                    numbers instead of
482             integers for easier computations.
483             :param state
484             :return discrete_state
485         """
486         min_states = self.env.observation_space.low
487         discrete_state = (state - min_states) * np.array
                ([10, 100])
488         return np.round(discrete_state, 0).astype(int)
489
490     def terminal_output(self, i):
491         # Terminal Output for stats of each epoch
492         avg_reward = sum(self.epoch_rewards[-2:]) / len(self
                .epoch_rewards[-2:])
493         self.epoch_rewards_table['ep'].append(i)
494         self.epoch_rewards_table['avg'].append(avg_reward)
495         self.epoch_rewards_table['min'].append(min(self.
                epoch_rewards))
496         self.epoch_rewards_table['max'].append(max(self.
                epoch_rewards))
497
498         print(f"Epoch - {i}\t| avg: {avg_reward:.2f}\t| min:
                {min(self.epoch_rewards[-1:]):.2f}"
499             f"\t| max: {max(self.epoch_rewards[-1:]):.2f}"
                    )
500
501 '''
502 Author: Caleb Ehrisman
503 Course- Advanced AI CSC-549
```

```python
504  Assignment - Programming Assignment #3 - Mountain Car
505  This file contains the code to implement the SARSA(lambda)
         algorithm.
506  All functions needed by solely the agent are included as
         member functions of class Agent
507  '''
508  import numpy as np
509  import itertools
510
511
512  class FourierBasis:
513      def __init__(self, state_space, order):
514          self.order = order
515          self.state_dim = state_space
516          self.order = [order]*self.state_dim
517          self.coeff = self.coefficients()
518
519      def coefficients(self):
520          """
521              FourierBasis.coefficients creates the coeffs for
                     the FourierBasis
522
523              :return np.array(coeff)
524          """
525          coeff = [np.zeros([self.state_dim])]
526
527          for i in range(0, self.state_dim):
528              for c in range(0, self.order[i]):
529                  v = np.zeros(self.state_dim)
530                  v[i] = c + 1
531                  coeff.append(v)
532          return np.array(coeff)
533
534      def get_features(self, state):
535          """
536              FourierBasis.get_features gets the feature
                     vector. Usually noted as x in SARSA(LAMBDA)
537
538              :param state
539
540              :return feature_vector
541          """
```

```
542             return np.cos(np.pi * np.dot(self.coeff, state))
```