

Project 3: Mountain Car

Sherwyn Braganza

December 3, 2022

Summary of Contents

- The Aim of the Project
- The Approach and Algorithm
- The Results
- Appendices

1 Aim of the Project

The aim of this project was to implement the famed 'Mountain Car' problem (Sutton and Barto, 1998) using SARSA(λ) with linear function approximation using Fourier Basis Functions.

The Mountain Car Problem involves a car in a Sinusoidal Valley who's aim is to reach the peak ([Figure 1](#)). The action space of the car consists of three actions - Left, Stationary and Right. The state space is two dimensional - $[x, \dot{x}]$ where x is the position and \dot{x} is the velocity of the body. Every action taken that doesn't lead to a terminal state (goal is reached), the reward is -1.

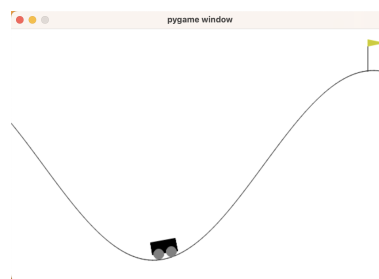


Figure 1: Mountain Car Environment

2 The Approach and Algorithm

2.1 Learner

In my implementation, I used SARSA(λ) as my learning agent. It exactly follows the algorithm shown in [Figure 2](#).

```

Input: a feature function  $\mathbf{x} : \mathcal{S}^+ \times \mathcal{A} \rightarrow \mathbb{R}^d$  such that  $\mathbf{x}(\text{terminal}, \cdot) = \mathbf{0}$ 
Input: a policy  $\pi$  (if estimating  $q_\pi$ )
Algorithm parameters: step size  $\alpha > 0$ , trace decay rate  $\lambda \in [0, 1]$ , small  $\epsilon > 0$ 
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  $\epsilon$ -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  $\epsilon$ -greedy according to  $\hat{q}(S', \cdot, \mathbf{w})$ 
     $\mathbf{x}' \leftarrow \mathbf{x}(S', A')$ 
     $Q \leftarrow \mathbf{w}^\top \mathbf{x}$ 
     $Q' \leftarrow \mathbf{w}^\top \mathbf{x}'$ 
     $\delta \leftarrow R + \gamma Q' - Q$ 
     $\mathbf{z} \leftarrow \gamma \lambda \mathbf{z} + (1 - \alpha \gamma \lambda \mathbf{z}^\top \mathbf{x}) \mathbf{x}$ 
     $\mathbf{w} \leftarrow \mathbf{w} + \alpha(\delta + Q - Q_{old})\mathbf{z} - \alpha(Q - Q_{old})\mathbf{x}$ 
     $Q_{old} \leftarrow Q$ 
     $\mathbf{x} \leftarrow \mathbf{x}'$ 
     $A \leftarrow A'$ 
  until  $S'$  is terminal

```

Figure 2: Online SARSA (λ)

\mathbf{w} is your weight vector and \mathbf{x} is the function approximation (in this case, a Fourier function approximation) for a given state and action.

α is the learning rate, λ is the trace decay rate and ϵ is a control factor that tells the algorithm whether it should follow the policy or take a exploratory action.

The only change that was made to this algorithm was that the episode would prematurely end if it took longer than 400 steps to reach the goal. This was done as a precautionary measure to mitigate penalty for starting in a bad state. The code that implements this algorithm can be found in [Appendix I](#).

2.2 Function Approximator

I used the method of Fourier Function Approximator described in [Value Function Approximation in Reinforcement Learning using the Fourier Basis by Konidaris et. al](#) to predict the Q values for a given state-action pair. The function took in two instantiation variables - the number of states and order that corresponded to the order of the Fourier basis function approximation. Using this it created a $(M + 1)^2 \times N$ matrix (\mathbf{C}), where N is the number of states and M is the order of the Fourier approximation. The **getFourierBasisApprox** function in my code used [Equation 1](#) to get ϕ or w or the Fourier basis which was then dotted with the weight vector in the SARSA algorithm to get the Q for a given state-action.

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

The code for this function can be found in [Appendix II](#)

3 Results

The program has two running modes:

1. Visual Mode - `python3 agent.py`
2. Analysis Mode - `python3 agent.py analysis`

The **Visual Mode** brings up a visual display of the algorithm learning. The **Analysis Mode** on the other hand doesn't have any GUI output. It however runs the learner using Fourier Approximations of base 3, 5 and 7. It performs 50 runs consisting of 250 episodes each and records the results. The data received for each basis function (3, 5 and 7) is averaged over the runs and then plotted against the other basis functions (Figure 3 shows this - the steps vs episodes plot).

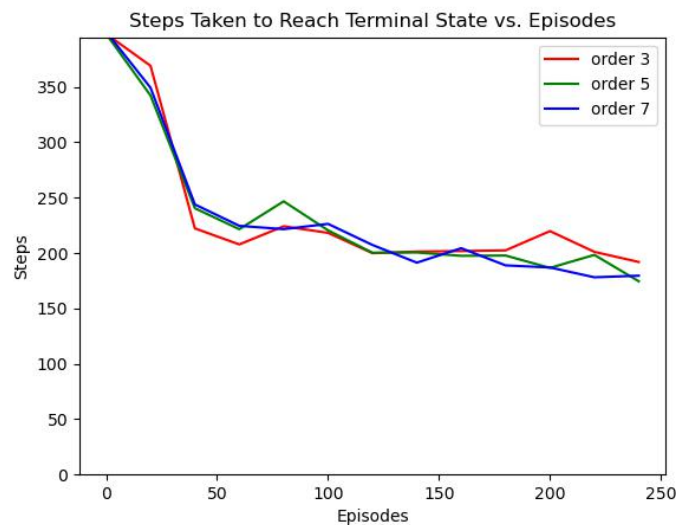
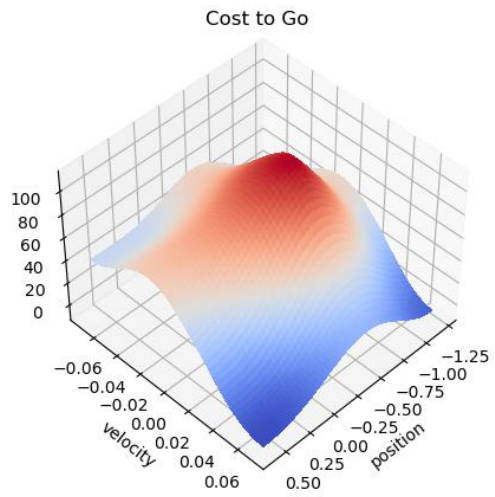
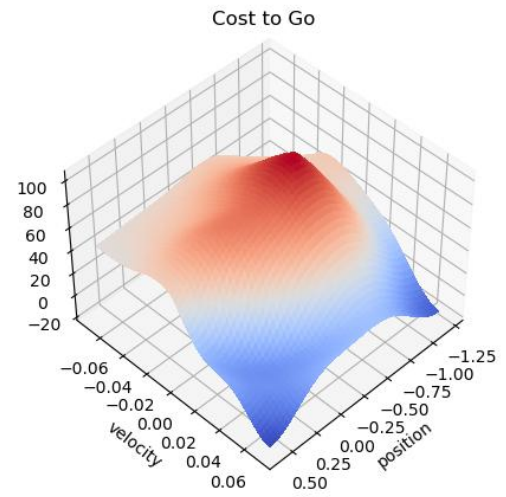


Figure 3: Steps vs Episodes in reaching terminal state

In addition to saving a copy of the previous plot, the program also prints the 'Cost to Go' plots for each Fourier Basis Approximation from the best trained models created for each basis. The plots Figure 4, Figure 5 and Figure 6 show the results obtained from one such run.



(a) Cost to Go - Order 3



(b) Cost to Go - Order 5

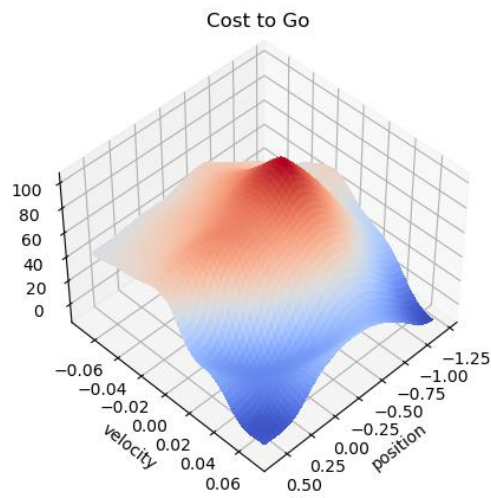


Figure 5: Cost to Go - Order 7

3.1 Questions

1) Show learning curves for order 3, 5, and 7 Fourier bases, for fixed setting of α and ϵ , and $\gamma = 1$, $\lambda = 0.9$.

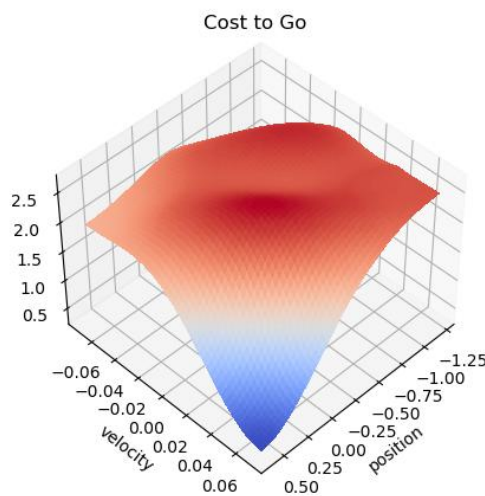
Figure 3 shows the plot obtained for the learning curves.

2) Create a surface plot of the value function (the negative of the value function) of the learned policies after 1, 000 episodes, for the above orders. (Hint: Your plot should look like the one in Sutton and Barto, but smoother.)

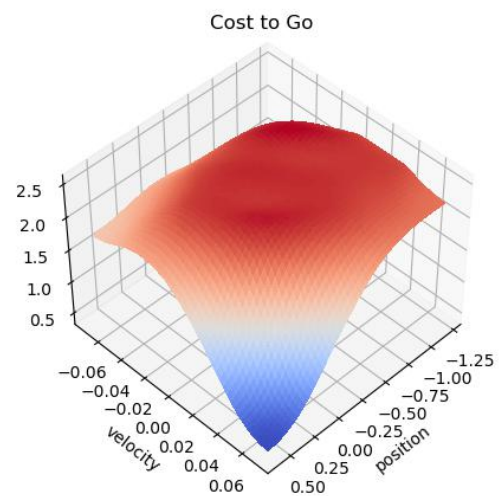
The plots Figure 4, Figure 5 and Figure 6 show this.

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

γ is the discount factor, which indicates how much the reward of a future state affects the current state. For a situation with a large number of states and a gamma less than 1, the effects of a profitable state in the future would die out really fast in previous states. This would make training slower when starting in bad initial states and would change the 'Cost to Go' plots drastically. The figures below show the Mountain Car Experiment with a $\gamma = 0.6$.



(a) Cost to Go - Order 3



(b) Cost to Go - Order 5

Consider the situation of having a path with a loop in it (Figure 7 below). In one scenario, the car travels through the loop to reach the right end and in the other, it bypasses the loop.

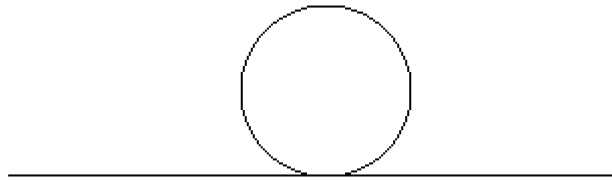


Figure 7: Track with Loop

If the reward for each step was changed to 0 and the final state 1, both episodes would be equally weighted even though one takes significantly more steps. Moreover, the episode that takes the loop would never end out of the loop. In this case a γ of 1 or less than that would not change anything

4 Appendix I

```

1 import copy
2 import numpy as np
3 import math
4 import sys
5 import random
6 import gym
7 from matplotlib import pyplot as plt
8 from function_approximator import FourierBasis
9 from matplotlib import cm
10
11
12 class SarsaLambdaLinear:
13     """
14     Implements SARSA Lamda with Linear Function Approximation
15
16     Attributes:
17         actions - The size of the action space
18         alpha - The learning rate
19         gamma - The discount factor
20         lamb - Lambda: the trace decay factor
21         epsilon - Factor that decides whether to follow policy or explore
22     """
23     def __init__(self, function_approximator, actions: int, alpha: float = 0.1,
24                  gamma: float = 0.99, lamb: float = 0.9, epsilon: float = 0.05, initial_valfunc: float = 0.0):
25         self.actions = actions
26         self.alpha = alpha
27         self.gamma = gamma
28         self.lamb = lamb
29         self.epsilon = epsilon
30         self.initial_valfunc = initial_valfunc
31         self.function_approximator = None
32
33         # if we don't have an approximation for each action, create deepcopies for each
34         if type(function_approximator) != list or len(function_approximator) < self.actions:
35             temp = []
36             for x in range(self.actions):
37                 temp.append(copy.deepcopy(function_approximator))
38
39             self.function_approximator = temp
40
41         if initial_valfunc == 0.0:
42             # initialize the 'theta' array corresponding to each action
43             self.weights = np.zeros([self.function_approximator[0].getShape()[0], self.actions])
44         else:
45             self.weights = np.ones([self.function_approximator[0].getShape()[0], self.actions]) * initial_valfunc
46
47         # initialize the weights for the trace update.
48         self.lambda_weights = np.zeros(self.weights.shape)
49
50     def traceClear(self):
51         """
52         Clear the weights of the trace vector
53         """
54         self.lambda_weights = np.zeros(self.weights.shape)
55
56     def makeOnPolicy(self):
57         """
58         Makes the policy greedy
59         """
60         self.epsilon = 0
61
62     def getStateActionVal(self, state, action):
63         """
64         Returns the Q value of the given state-action pair.
65         """
66         return np.dot(self.weights[:, action], self.function_approximator[action].getFourierBasisApprox(state))
67
68     def getMaxStateActionVal(self, state):
69         """
70         Checks all Q values in
71         """
72         best = float('-inf')
73         best_a = float('-inf')
74         for a in range(0, self.actions):
75             qval = self.getStateActionVal(state, a)
76             if qval > best:
77                 best = qval
78                 best_a = a
79
80         return best, best_a
81
82     def next_move(self, state):
83         """
84         Stochastically returns an epsilon-greedy action from the current state.
85         """
86
87         if np.random.random() <= self.epsilon:
88             return np.random.randint(0, self.actions)

```

```

89
90 # Build a list of actions evaluating to max_a Q(s, a)
91 best = float("-inf")
92 best_actions = []
93
94 for a in range(self.actions):
95     thisq = self.getStateActionVal(state, a)
96
97     if abs(thisq - best) < 0.001:
98         best_actions.append(a)
99     elif thisq > best:
100         best = thisq
101         best_actions = [a]
102
103 if len(best_actions) == 0 or math.isinf(best):
104     print("SarsaLambdaLinearFA: function approximator has diverged to infinity.", file=sys.stderr)
105     return np.random.randint(0, self.actions)
106
107 # Select randomly among best-valued actions
108 return random.choice(best_actions)
109
110 def update(self, state, action, reward, next_state, next_action=None, terminal=False) -> float:
111     """
112     Runs a Sarsa update, given a transition.
113     If no action is provided, it assumes an E-Greedy policy and finds
114     an action that maximizes the Q value.
115
116     Parameters
117     -----
118     state - the state at time t
119     action - the action to be taken to reach s+1
120     reward - The reward received
121     next_state - the state at t+1
122     next_action - the action at t+1, if not present it is calculated
123     terminal - if the next state is the terminal state.
124
125     Returns
126     -----
127     delta - The TD error.
128
129     """
130
131     # Compute TD error
132     delta = reward - self.getStateActionVal(state, action)
133
134     # Only include s' if it is not a terminal state.
135     if not terminal:
136         if next_action is not None:
137             delta += self.gamma * self.getStateActionVal(next_state, next_action)
138         else:
139             # adopt an exploration action
140             (next_Q, next_action) = self.getMaxStateActionVal(next_state)
141             delta += self.gamma * self.getMaxStateActionVal(next_Q)
142
143     # Compute the basis functions for state s, action a.
144     eval_f_action = self.function_approximator[action].getFourierBasisApprox(state)
145
146     for each_a in range(0, self.actions):
147
148         # Trace Update
149         self.lambda_weights[:, each_a] += self.gamma * self.lamb
150         if each_a == action:
151             self.lambda_weights[:, each_a] += eval_f_action
152
153         # Weight Update
154         self.weights[:, each_a] += self.alpha * \
155             delta * \
156             np.multiply(self.function_approximator[each_a].getGradientFactors(),
157                         self.lambda_weights[:, each_a])
158
159     # Return the TD error, which may be informative.
160     return delta
161
162 def fourierBasis(env, samples: int = 10, episodes: int = 10, order: int = 3):
163     gamma = 1.0
164     run_data = np.zeros((samples, episodes, 2))
165
166     state_dim = env.observation_space.shape[0]
167     actions = env.action_space.n
168     u_state = env.observation_space.high
169     l_state = env.observation_space.low
170     d_state = u_state - l_state
171     best_learner = None
172     best_sum = float("-inf")
173
174     for sample in range(0, samples):
175
176         fb = FourierBasis(order=order, dimensions=state_dim)
177         learner = SarsaLambdaLinear(fb, actions=actions, gamma=gamma, lamb=0.95, epsilon=0.05, alpha=0.001)
178
179         for episode in range(0, episodes):
180             steps = 0

```



```

181     # converge to pure on-policy for last 10 episodes
182     if episode >= 0.8 * episodes:
183         learner.makeOnPolicy()
184
185     learner.traceClear()
186     s = (env.reset() - l.state) / d.state
187     a = learner.next.move(s)
188
189     done = False
190     nsteps = 0
191     sum_r = 0.0
192
193     while not done:
194         sp, r, done, info = env.step(a)
195         steps += 1
196         sp = (sp - l.state) / d.state
197         term = (done and not (info.get('TimeLimit.truncated', False)))
198         ap = learner.next.move(sp)
199
200         learner.update(s, a, r, sp, ap, terminal=term)
201         s = sp
202         a = ap
203         sum_r += r * pow(gamma, nsteps)
204         steps += 1
205
206     run_data[sample, episode, :] = np.asarray([sum_r, steps])
207     best_learner = learner if sum_r >= best_sum else best_learner
208
209     # print('Run ' + str(sample + 1) + ", ep. " + str(episode + 1) + " return: " + str(sum_r) +
210     #       ", # steps: " + str(steps))
211
212     env.close()
213
214     return run_data, best_learner
215
216 def runAnalysis():
217     run_data = []
218     learner = []
219     basis = [3, 5, 7]
220     colors = ['red', 'green', 'blue']
221     episodes = 250
222     samples = 50
223
224     env = gym.make('MountainCar-v0')
225
226     u_state = env.observation_space.high
227     l_state = env.observation_space.low
228     d_state = u_state - l_state
229
230     for i in basis:
231         temp_data, temp_learner = fourierBasis(env, samples, episodes, i)
232         run_data.append(temp_data)
233         learner.append(temp_learner)
234
235     run_data_mean = []
236     # run_data_std_dev = []
237     data_range = range(0, episodes)
238
239     for i in range(len(basis)):
240         run_data_mean.append(np.mean(run_data[i], axis=0))
241         # run_data_std_dev.append(np.std(run_data[i], axis=0))
242         plt.plot(data_range[0::20], run_data_mean[i][0::20, 1], c=colors[i], label='order ' + str(basis[i]))
243
244     # plotting steps taken to reach term vs. episodes
245     plt.xlabel('Episodes')
246     plt.ylabel('Steps')
247     plt.ylim(0, 395)
248     plt.title('Steps Taken to Reach Terminal State vs. Episodes')
249     plt.legend()
250     plt.savefig('steps-vs-episodes.jpg')
251     plt.close()
252
253     # value function surface plots
254     x_bounds, y_bounds = [l_state[0], u_state[0]], [l_state[1], u_state[1]]
255     xs, ys = np.meshgrid(np.linspace(x_bounds[0], x_bounds[1], 100),
256                          np.linspace(y_bounds[0], y_bounds[1], 100))
257     zs = np.zeros(xs.shape)
258
259     for base in range(len(basis)):
260         for i in range(0, zs.shape[0]):
261             for j in range(0, zs.shape[1]):
262                 s = [(xs[i, j] - l_state[0]) / d_state[0], (ys[i, j] - l_state[1]) / d_state[1]]
263                 (zq, _) = learner[base].getMaxStateActionVal(s)
264                 zs[i, j] = -1.0 * zq
265
266     fig = plt.figure()
267     ax = fig.gca(projection='3d')
268     ax.plot_surface(xs, ys, zs, cmap=cm.get_cmap("coolwarm"), linewidth=0, antialiased=False)
269     ax.view_init(elev=45, azim=45)
270     ax.set_xlabel('position')
271     ax.set_ylabel('velocity')
272     ax.view_init(elev=45, azim=45)

```

```
273     ax.set_title('Cost to Go')
274     fig.savefig('Cost to Go - Order ' + str(basis[base]) + '.jpeg')
275     plt.close()
276
277 def runVisualDisplay():
278     episodes = 250
279     samples = 10
280     env = gym.make('MountainCar-v0', render_mode='human').env
281     results, learner = fourierBasis(env, samples=samples, episodes=episodes, order=5)
282
283
284 if __name__ == '__main__':
285     if len(sys.argv) > 1:
286         if sys.argv[1] == 'analysis':
287             runAnalysis()
288         else:
289             runVisualDisplay()
```

5 Appendix II

```

1 class FourierBasis:
2     def __init__(self, order: int, dimensions: int):
3         # Instance variables
4         self.coefficients = np.array([])
5         self.gradient_factors = np.array([])
6         self.dimensions = dimensions
7         self.order = [order] * self.dimensions
8
9         # create empty container for coefficient array
10        prods = [range(0, o + 1) for o in self.order]
11        coeffs = [v for v in itertools.product(*prods)]
12        self.coefficients = np.array(coeffs)
13
14        with np.errstate(divide='ignore', invalid='ignore'):
15            self.gradient_factors = 1.0 / np.linalg.norm(self.coefficients, ord=2, axis=1)
16            self.gradient_factors[0] = 1.0 # Overwrite division by zero for function with all-zero coefficients.
17
18    def getFourierBasisApprox(self, state_vector: np.ndarray):
19        """
20        Computes basis function values at a given state.
21        """
22
23        # Bounds check state vector
24        if np.min(state_vector) < 0.0 or np.max(state_vector) > 1.0:
25            print('Fourier Basis: Given State Vector ({} ) not in range [0.0, 1.0]'.format(state_vector),
26                  file=sys.stderr)
27
28        # Compute the Fourier Basis feature values
29        return np.cos(np.pi * np.dot(self.coefficients, state_vector))
30
31    def getShape(self):
32        return self.coefficients.shape
33
34    def getGradientFactors(self):
35        return self.gradient_factors
36
37    def getGradientFactor(self, function_no):
38        return self.gradient_factors[function_no]
39
40    def length(self):
41        """Return the number of basis functions."""
42        return self.coefficients.shape[0]

```