

```

1  """
2  CS 229 Machine Learning
3  Question: Reinforcement Learning - The Inverted Pendulum
4  """
5  from __future__ import division, print_function
6  from env import CartPole, Physics
7  import matplotlib.pyplot as plt
8  import numpy as np
9  from scipy.signal import lfilter
10
11  """
12  Parts of the code (cart and pole dynamics, and the state
13  discretization) are inspired from code available at the RL repository
14  http://www-anw.cs.umass.edu/rlr/domains.html
15
16  Briefly, the cart-pole system is described in `cart_pole.py`. The main
17  simulation loop in this file calls the `simulate()` function for
18  simulating the pole dynamics, `get_state()` for discretizing the
19  otherwise continuous state space in discrete states, and `show_cart()`
20  for display.
21
22  Some useful parameters are listed below:
23
24  `NUM_STATES`: Number of states in the discretized state space
25  You must assume that states are numbered 0 through `NUM_STATES` - 1. The
26  state numbered `NUM_STATES` - 1 (the last one) is a special state that
27  marks the state when the pole has been judged to have fallen (or when
28  the cart is out of bounds). However, you should NOT treat this state
29  any differently in your code. Any distinctions you need to make between
30  states should come automatically from your learning algorithm.
31
32  After each simulation cycle, you are supposed to update the transition
33  counts and rewards observed. However, you should not change either
34  your value function or the transition probability matrix at each
35  cycle.
36
37  Whenever the pole falls, a section of your code below will be
38  executed. At this point, you must use the transition counts and reward
39  observations that you have gathered to generate a new model for the MDP
40  (i.e. transition probabilities and state rewards). After that, you
41  must use value iteration to get the optimal value function for this MDP
42  model.
43
44  `TOLERANCE`: Controls the convergence criteria for each value iteration
45  run. In value iteration, you can assume convergence when the maximum
46  absolute change in the value function at any state in an iteration
47  becomes lower than `TOLERANCE`.
48
49  You need to write code that chooses the best action according
50  to your current value function, and the current model of the MDP. The
51  action must be either 0 or 1 (corresponding to possible directions of
52  pushing the cart)
53
54  Finally, we assume that the simulation has converged when
55  `NO_LEARNING_THRESHOLD` consecutive value function computations all
56  converged within one value function iteration. Intuitively, it seems
57  like there will be little learning after this, so we end the simulation
58  here, and say the overall algorithm has converged.
59
60
61  Learning curves can be generated by calling a code snippet at the end
62  (it assumes that the learning was just executed, and the array
63  `time_steps_to_failure` that records the time for which the pole was
64  balanced before each failure is in memory). `num_failures` is a variable
65  that stores the number of failures (pole drops / cart out of bounds)
66  till now.
67
68  Other parameters in the code are described below:
69
70  `GAMMA`: Discount factor to be used
71
72  The following parameters control the simulation display; you dont

```



```

73 really need to know about them:
74
75 `pause_time`: Controls the pause between successive frames of the
76 display. Higher values make your simulation slower.
77 `min_trial_length_to_start_display`: Allows you to start the display only
78 after the pole has been successfully balanced for at least this many
79 trials. Setting this to zero starts the display immediately. Choosing a
80 reasonably high value (around 100) can allow you to rush through the
81 initial learning quickly, and start the display only after the
82 performance is reasonable.
83 """
84
85 def initialize_mdp_data(num_states):
86     """
87     Return a variable that contains all the parameters/state you need for your MDP.
88     Feel free to use whatever data type is most convenient for you (custom classes, tuples, dicts, etc)
89
90     Assume that no transitions or rewards have been observed.
91     Initialize the value function array to small random values (0 to 0.10, say).
92     Initialize the transition probabilities uniformly (ie, probability of
93     transitioning for state x to state y using action a is exactly
94     1/num_states).
95     Initialize all state rewards to zero.
96
97     Args:
98         num_states: The number of states
99
100     Returns: The initial MDP parameters
101     """
102     transition_counts = np.zeros((num_states, num_states, 2))
103     transition_probs = np.ones((num_states, num_states, 2)) / num_states
104     #Index zero is count of rewards being -1 , index 1 is count of total num state is reached
105     reward_counts = np.zeros((num_states, 2))
106     reward = np.zeros(num_states)
107     value = np.random.rand(num_states) * 0.1
108
109     return {
110         'transition_counts': transition_counts,
111         'transition_probs': transition_probs,
112         'reward_counts': reward_counts,
113         'reward': reward,
114         'value': value,
115         'num_states': num_states,
116     }
117
118 def choose_action(state, mdp_data):
119     """
120     Choose the next action (0 or 1) that is optimal according to your current
121     mdp_data. When there is no optimal action, return a random action.
122
123     Args:
124         state: The current state in the MDP
125         mdp_data: The parameters for your MDP. See initialize_mdp_data.
126
127     Returns:
128         0 or 1 that is optimal according to your current MDP
129     """
130
131     # *** START CODE HERE ***
132     score1 = mdp_data['transition_probs'][state, :, 0].dot(mdp_data['value'])
133     score2 = mdp_data['transition_probs'][state, :, 1].dot(mdp_data['value'])
134
135     if score1 > score2:
136         action = 0
137     elif score2 > score1:
138         action = 1
139     else:
140         action = 0 if np.random.uniform() < 0.5 else 1
141
142     return action
143     # *** END CODE HERE ***
144

```

```

145 def update_mdp_transition_counts_reward_counts(mdp_data, state, action, new_state, reward):
146     """
147     Update the transition count and reward count information in your mdp_data.
148     Do not change the other MDP parameters (those get changed later).
149
150     Record the number of times `state, action, new_state` occurs.
151     Record the rewards for every `new_state`
152     (since rewards are -1 or 0, you just need to record number of times reward -1 is seen in 'reward_counts' index new_state,0)
153     Record the number of time `new_state` was reached (in 'reward_counts' index new_state,1)
154
155     Args:
156         mdp_data: The parameters of your MDP. See initialize_mdp_data.
157         state: The state that was observed at the start.
158         action: The action you performed.
159         new_state: The state after your action.
160         reward: The reward after your action (i.e. reward corresponding to new_state).
161
162     Returns:
163         Nothing
164     """
165
166     # *** START CODE HERE ***
167     mdp_data['transition_counts'][state, new_state, action] += 1
168     if reward == -1:
169         mdp_data['reward_counts'][new_state, 0] += 1
170     mdp_data['reward_counts'][new_state, 1] += 1
171     # *** END CODE HERE ***
172
173     # This function does not return anything
174     return
175
176 def update_mdp_transition_probs_reward(mdp_data):
177     """
178     Update the estimated transition probabilities and reward values in your MDP.
179
180     Make sure you account for the case when a state-action pair has never
181     been tried before, or the state has never been visited before. In that
182     case, you must not change that component (and thus keep it at the
183     initialized uniform distribution).
184
185     Args:
186         mdp_data: The data for your MDP. See initialize_mdp_data.
187
188     Returns:
189         Nothing
190     """
191
192     # *** START CODE HERE ***
193     for a in [0, 1]:
194         for s in range(mdp_data['num_states']):
195             total_num_transitions = np.sum(mdp_data['transition_counts'][s, :, a])
196             if total_num_transitions > 0:
197                 mdp_data['transition_probs'][s, :, a] = (
198                     mdp_data['transition_counts'][s, :, a] / total_num_transitions
199                 )
200
201     for s in range(mdp_data['num_states']):
202         if mdp_data['reward_counts'][s, 1] > 0:
203             mdp_data['reward'][s] = -mdp_data['reward_counts'][s, 0] / mdp_data['reward_counts'][s, 1]
204
205     # *** END CODE HERE ***
206
207     # This function does not return anything
208     return
209
210 def update_mdp_value(mdp_data, tolerance, gamma):
211     """
212     Update the estimated values in your MDP.
213
214
215

```



```

216 Perform value iteration using the new estimated model for the MDP.
217 The convergence criterion should be based on 'TOLERANCE' as described
218 at the top of the file.
219
220 Return true if it converges within one iteration.
221
222 Args:
223     mdp_data: The data for your MDP. See initialize_mdp_data.
224     tolerance: The tolerance to use for the convergence criterion.
225     gamma: Your discount factor.
226
227 Returns:
228     True if the value iteration converged in one iteration
229
230 """
231
232 # *** START CODE HERE ***
233 iterations = 0
234
235 while True:
236     new_value = np.zeros(mdp_data['num_states'])
237
238     iterations = iterations + 1
239     for s in range(mdp_data['num_states']):
240         value1 = mdp_data['transition_probs'][s, :, 0].dot(mdp_data['value'])
241         value2 = mdp_data['transition_probs'][s, :, 1].dot(mdp_data['value'])
242
243         new_value[s] = max(value1, value2)
244
245     new_value = mdp_data['reward'] + gamma * new_value
246
247     max_diff = float('-Inf')
248     for s in range(mdp_data['num_states']):
249         if abs(new_value[s] - mdp_data['value'][s]) > max_diff:
250             max_diff = abs(new_value[s] - mdp_data['value'][s])
251
252     mdp_data['value'] = new_value
253
254     if max_diff < tolerance:
255         break
256
257 return iterations == 1
258
259 # *** END CODE HERE ***
260
261 def main(plot=True):
262     # Seed the randomness of the simulation so this outputs the same thing each time
263     np.random.seed(0)
264
265     # Simulation parameters
266     pause_time = 0.0001
267     min_trial_length_to_start_display = 100
268     display_started = min_trial_length_to_start_display == 0
269
270     NUM_STATES = 163
271     GAMMA = 0.995
272     TOLERANCE = 0.01
273     NO_LEARNING_THRESHOLD = 20
274
275     # Time cycle of the simulation
276     time = 0
277
278     # These variables perform bookkeeping (how many cycles was the pole
279     # balanced for before it fell). Useful for plotting learning curves.
280     time_steps_to_failure = []
281     num_failures = 0
282     time_at_start_of_current_trial = 0
283
284     # You should reach convergence well before this
285     max_failures = 500
286
287     # Initialize a cart pole

```

```

288 cart_pole = CartPole(Physics())
289
290 # Starting `state_tuple` is (0, 0, 0, 0)
291 # x, x_dot, theta, theta_dot represents the actual continuous state vector
292 x, x_dot, theta, theta_dot = 0.0, 0.0, 0.0, 0.0
293 state_tuple = (x, x_dot, theta, theta_dot)
294
295 # `state` is the number given to this state, you only need to consider
296 # this representation of the state
297 state = cart_pole.get_state(state_tuple)
298 # if min_trial_length_to_start_display == 0 or display_started == 1:
299 #     cart_pole.show_cart(state_tuple, pause_time)
300
301 mdp_data = initialize_mdp_data(NUM_STATES)
302
303 # This is the criterion to end the simulation.
304 # You should change it to terminate when the previous
305 # 'NO_LEARNING_THRESHOLD' consecutive value function computations all
306 # converged within one value function iteration. Intuitively, it seems
307 # like there will be little learning after this, so end the simulation
308 # here, and say the overall algorithm has converged.
309
310 consecutive_no_learning_trials = 0
311 while consecutive_no_learning_trials < NO_LEARNING_THRESHOLD:
312
313     action = choose_action(state, mdp_data)
314
315     # Get the next state by simulating the dynamics
316     state_tuple = cart_pole.simulate(action, state_tuple)
317     # x, x_dot, theta, theta_dot = state_tuple
318
319     # Increment simulation time
320     time = time + 1
321
322     # Get the state number corresponding to new state vector
323     new_state = cart_pole.get_state(state_tuple)
324     # if display_started == 1:
325     #     cart_pole.show_cart(state_tuple, pause_time)
326
327     # reward function to use - do not change this!
328     if new_state == NUM_STATES - 1:
329         R = -1
330     else:
331         R = 0
332
333     update_mdp_transition_counts_reward_counts(mdp_data, state, action, new_state, R)
334
335     # Recompute MDP model whenever pole falls
336     # Compute the value function V for the new model
337     if new_state == NUM_STATES - 1:
338
339         update_mdp_transition_probs_reward(mdp_data)
340
341         converged_in_one_iteration = update_mdp_value(mdp_data, TOLERANCE, GAMMA)
342
343         if converged_in_one_iteration:
344             consecutive_no_learning_trials = consecutive_no_learning_trials + 1
345         else:
346             consecutive_no_learning_trials = 0
347
348     # Do NOT change this code: Controls the simulation, and handles the case
349     # when the pole fell and the state must be reinitialized.
350     if new_state == NUM_STATES - 1:
351         num_failures += 1
352         if num_failures >= max_failures:
353             break
354         print('[INFO] Failure number {}'.format(num_failures))
355         time_steps_to_failure.append(time - time_at_start_of_current_trial)
356         # time_steps_to_failure[num_failures] = time - time_at_start_of_current_trial
357         time_at_start_of_current_trial = time
358
359     if time_steps_to_failure[num_failures - 1] > min_trial_length_to_start_display:

```

```

360         display_started = 1
361
362         # Reinitialize state
363         # x = 0.0
364         x = -1.1 + np.random.uniform() * 2.2
365         x_dot, theta, theta_dot = 0.0, 0.0, 0.0
366         state_tuple = (x, x_dot, theta, theta_dot)
367         state = cart_pole.get_state(state_tuple)
368     else:
369         state = new_state
370
371     if plot:
372         # plot the learning curve (time balanced vs. trial)
373         log_tstf = np.log(np.array(time_steps_to_failure))
374         plt.plot(np.arange(len(time_steps_to_failure)), log_tstf, 'k')
375         window = 30
376         w = np.array([1/window for _ in range(window)])
377         weights = lfilter(w, 1, log_tstf)
378         x = np.arange(window//2, len(log_tstf) - window//2)
379         plt.plot(x, weights[window:len(log_tstf)], 'r--')
380         plt.xlabel('Num failures')
381         plt.ylabel('Log of num steps to failure')
382         plt.savefig('./control.pdf')
383
384     return np.array(time_steps_to_failure)
385
386 if __name__ == '__main__':
387     main()

```