

# Project 3: Approximation Methods and Policy Gradients

August Jonasson

March 10, 2025

## Part 1: gradients of update steps in SGD

In this first part we consider the parameter update steps in approximation methods (Sutton 2018, ch. 9) and policy gradient methods (ch. 13). More specifically, the tasks will consist of deriving the gradients needed in order to perform these update steps, given some model assumptions.

### Task 1

In approximation methods of the state-value  $v_\pi$ , we use SGD with the update rule

$$\theta_{t+1} = \theta_t + \alpha \delta_t \nabla_\theta \hat{v}(s, \theta_t), \quad (1)$$

where  $\alpha$  is the learning rate,  $\delta_t = v_\pi(s) - \hat{v}(s, \theta_t)$  is the estimation error at time  $t$ , and  $\nabla_\theta$  denotes the gradient w.r.t. the parameter vector  $\theta$ .

Suppose we have a linear approximator  $\hat{v}(s, \theta) = \theta^T x(s)$ , where  $x(s)$  denotes some feature of state  $s$ , and show that  $\theta$  updates according to

$$\theta_{t+1} = \theta_t + \alpha \delta_t x(s).$$

Using (1), all we have to show is that  $\nabla_\theta \hat{v}(s, \theta) = x(s)$ , i.e. that

$$\frac{\partial(\hat{v}(s, \theta))}{\partial \theta} = x(s).$$

Let  $\theta, x(s) \in \mathbb{R}^p$ , then (using the short-hand  $\hat{v} = \hat{v}(s, \theta)$ )

$$\frac{\partial \hat{v}}{\partial \theta} = \left[ \frac{\partial \hat{v}}{\partial \theta_1}, \dots, \frac{\partial \hat{v}}{\partial \theta_p} \right]^T,$$

with

$$\frac{\partial \hat{v}}{\partial \theta_i} = \frac{\partial(\theta^T x(s))}{\partial \theta_i} = \frac{\partial}{\partial \theta_i} \left( \sum_{j=1}^p \theta_j x(s)_j \right) = x(s)_i,$$

for  $1 \leq i \leq p$ . Thus

$$\frac{\partial \hat{v}}{\partial \theta} = [x(s)_1, \dots, x(s)_p]^T = x(s),$$

and we are done.

## Task 2

For policy gradient methods, we update the weights according to

$$\theta_{t+1} = \theta_t + \alpha \delta_t \frac{\nabla_{\theta} \pi(a|s; \theta_t)}{\pi(a|s; \theta_t)}. \quad (2)$$

(a)

Show that

$$\frac{\nabla_{\theta} \pi(a|s; \theta)}{\pi(a|s; \theta)} = \nabla_{\theta} \log \pi(a|s; \theta).$$

Starting from the right-hand side (with short-hand  $\pi = \pi(a|s; \theta)$ ), we have that

$$\nabla_{\theta} \log \pi = \left[ \frac{\partial \log \pi}{\partial \theta_1}, \dots, \frac{\partial \log \pi}{\partial \theta_p} \right]^T,$$

where

$$\frac{\partial \log \pi}{\partial \theta_i} = \frac{1}{\pi} \frac{\partial \pi}{\partial \theta_i}$$

for  $1 \leq i \leq p$ . Thus

$$\nabla_{\theta} \log \pi = \frac{1}{\pi} \left[ \frac{\partial \pi}{\partial \theta_1}, \dots, \frac{\partial \pi}{\partial \theta_p} \right] = \frac{\nabla_{\theta} \pi(a|s; \theta)}{\pi(a|s; \theta)}, \quad (3)$$

and we are done.

(b)

Suppose the policy is modeled using a softmax function

$$\pi(a|s; \theta) = \frac{e^{h(s,a,\theta)}}{\sum_b e^{h(s,b,\theta)}}, \quad (4)$$

with linear input  $h(s, a, \theta) = \theta^T x(s, a)$  for some state-action feature function  $x$ . Show that the policy update in (2) is given by

$$\theta_{t+1} = \theta_t + \alpha \delta_t \left( x(s, a) - \sum_b \pi(b|s; \theta_t) x(s, b) \right).$$

Using (2) and (3), all we have to show is that

$$\nabla_{\theta} \log \pi(a|s; \theta) = x(s, a) - \sum_b \pi(b|s; \theta) x(s, b).$$

Starting from the left-hand side, we first have that

$$\log \pi(a|s; \theta) = \theta^T x(s, a) - \log \left( \sum_b e^{\theta^T x(s, b)} \right),$$

from which using the results from Task 1, the partial derivatives w.r.t.  $\theta$  of the first term can be directly simplified to the sought after  $x(s, a)$ . The derivatives of the second term can be manipulated in the following way:

$$\begin{aligned}\frac{\partial}{\partial \theta} \log \left( \sum_b e^{\theta^T x(s, b)} \right) &= \frac{1}{\sum_b e^{\theta^T x(s, b)}} \frac{\partial}{\partial \theta} \left( \sum_b e^{\theta^T x(s, b)} \right) \\ &= \frac{1}{\sum_b e^{\theta^T x(s, b)}} \left( \sum_b \frac{\partial}{\partial \theta} e^{\theta^T x(s, b)} \right),\end{aligned}$$

and again, using the results from Task 1, the partial derivative within the sum reduces to  $x(s, b)e^{\theta^T x(s, b)}$ , for all  $b$ . Putting it all together, we now have

$$\begin{aligned}\nabla_{\theta} \log \pi(a|s; \theta) &= x(s, a) - \frac{1}{\sum_c e^{\theta^T x(s, c)}} \sum_b x(s, b) e^{\theta^T x(s, b)} \\ &= x(s, a) - \sum_b x(s, b) \frac{e^{\theta^T x(s, b)}}{\sum_c e^{\theta^T x(s, c)}} \\ &= x(s, a) - \sum_b x(s, b) \pi(b|s; \theta),\end{aligned}$$

and we are done.

## Part 2: approximation and policy gradient methods for grid-world with a monster

In this second part we are to implement approximation and policy gradient methods in code, for the gridworld with a monster that we worked with in the previous project. Refer back to Project 2 for a full description of system dynamics. We will now state all information that is new.

We set the grid size to  $N = 10$ , episode length  $T = 200$  and discount factor  $\gamma = 0.95$ . Furthermore, we define the features

$$\begin{aligned}x(s, a) &= [f_1, f_2, f_3, f_4]^T, \quad \text{where} \\ f_1 &= \frac{1}{\text{distance from agent to apple} + 1}, \\ f_2 &= \frac{1}{\text{distance from agent to monster} + 1}, \\ f_3 &= 1 \text{ if action takes agent closer to apple, else } 0 \\ f_4 &= 1 \text{ if action takes agent closer to monster, else } 0,\end{aligned}$$

with the Manhattan distance used as a measure of distance since we are operating on a gridworld. Furthermore, we use the linear approximators for state and action-value

$$\begin{aligned}\hat{q}(s, a, \theta) &= \theta^T x(s, a) \quad \text{and} \\ \hat{v}(s) &= \max_a \hat{q}(s, a, \theta),\end{aligned}$$

and for the policy gradient methods we model the policy as in (4) - using a softmax function with linear input.

## Task 1

Implement the following methods:

- Semi-gradient n-step SARSA for  $n = 1, 2, 3$ , with a suitable  $\epsilon$  (p. 247 in course book),
- REINFORCE (p. 328 in course book),
- REINFORCE with baseline (p. 330 in course book), and
- One-step Actor-Critic (p. 332 in course book).

See code for full implementation.

## Task 2

We are to plot the learning curves for the methods implemented in the previous task.

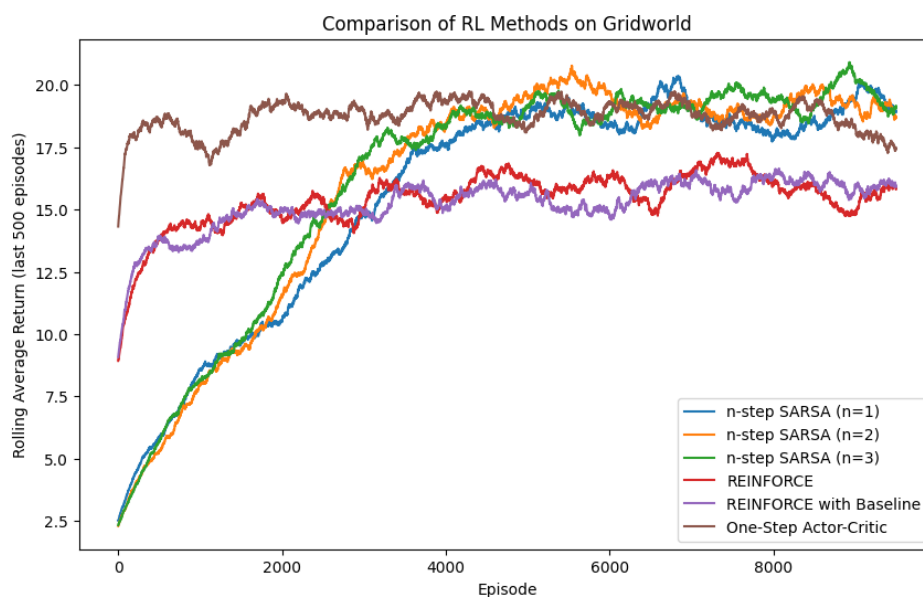


Figure 1: Rolling average (previous 500) of total reward per episode for the different approximation and policy gradient methods.

## Task 3

Discuss the following:

1. How does the learning rate  $\alpha$  affect the learning speed and convergence of the methods?

2. How do the methods above compare in terms of convergence speed?
3. In this example, does policy gradient methods perform better than the value-based methods? Why or why not?

(1)

Higher learning rate  $\alpha$  might lead to faster initial learning but risks continuously over-shooting the “true solution” by a lot, thus leading to oscillations in the learning curve. In general, a smaller  $\alpha$  is more safe in the sense that it is more stable and more precise, but comes at the price of slower learning. In most of the methods in Figure 1, we found that  $\alpha = 0.01$  could lead to unstable solutions (big oscillations), and thus, we opted for  $\alpha = 0.001$  in all methods.

(2)

Methods that use bootstrapping (like one-step actor-critic or n-step SARSA) generally converge faster than pure Monte Carlo methods (e.g. plain REINFORCE). This is because bootstrapping incorporates information from future estimates, which can lead to quicker adjustments of values and policies. So why do the n-step methods in Figure 1 converge the slowest? It is because we are using an exponentially decaying  $\epsilon$ . At episode 3000, the exploration is still at approximately 23 %, before it at around 6500 episodes drops below 5 %. With this in mind, it is not surprising that the reward has yet to stabilize after so many episodes. We are sacrificing convergence speed for a better policy - and a better policy we do indeed find. With a constant  $\epsilon$ , the n-step methods performed the worst among all methods. Now, they are arguably the best - as long as we can afford the extra time they take to train.

Among the remaining three episodes, it is difficult to say which one converges the fastest. Based on Figure 1, they all seem rather similar. We can however say that we expect the actor-critic and REINFORCE with baseline methods to converge faster than plain REINFORCE, due to the continual bootstrapping from the critic and baseline respectively. Maybe in this setting, the differences are marginal.

(3)

If we disregard the time it takes to train and focus only on which method achieves the best policy, then the value-based methods (n-step SARSA) seem to perform the best. This is most likely due to the fact that the gridworld system is a very tabular setting. Relatively speaking, there are not that many state-action pairs to consider, and we know that the policy gradient methods really shine the most when the action space is discrete, but the state space continuous.

In summary: gridworld is a low-dimensional well structured and discrete problem. In cases such as these, the simpler value-based methods can capture the necessary dynamics quickly without the additional complexity introduced by additional parameters of e.g. actor-critic or baseline methods. Even though the actor-critic (or REINFORCE with baseline) method may perform as well as the n-step methods, it is simply unnecessary to introduce the extra complexity that comes with policy gradient methods.

## Appendix 1: Python code

```
1 # %%
2 import numpy as np
3 import random
4 from matplotlib import pyplot as plt
5
6 # -----
7 # Environment Setup
8 # -----
9 N = 10 # grid size (N x N)
10 T = 200 # maximum time steps per episode
11 alpha = 0.001 # learning rate
12 gamma = 0.95 # discount factor
13 episodes = 10**4 # number of episodes per method
14
15 # Create grid coordinates
16 grid = [(r, c) for r in range(N) for c in range(N)]
17
18 # Define actions and their effects (using (row, col) convention)
19 actions = ["up", "right", "down", "left"]
20 actions_dict = {
21     "up": (-1, 0),
22     "right": (0, 1),
23     "down": (1, 0),
24     "left": (0, -1)
25 }
26
27 # Define rewards
28 rewards_dict = {
29     "collect_apple": 1,
30     "caught_by_monster": -1,
31     "empty": 0
32 }
33
34 # Parameters for epsilon decay (used in n-step SARSA)
35 epsilon_0 = 1.0
36 epsilon_min = 0.01
37 decay_rate = 0.0005
38
39 def get_epsilon(episode):
40     """Exponentially decaying epsilon."""
41     return epsilon_min + (epsilon_0 - epsilon_min) * np.exp(-decay_rate * episode)
42
43 def initial_positions():
44     """Return distinct starting positions for agent, monster, and apple."""
45     return tuple(random.sample(grid, 3))
46
47 def respawn_apple(agent_pos, monster_pos):
48     """Return a new apple position that is not occupied by agent or monster."""
49     available_positions = [pos for pos in grid if pos != agent_pos and pos != monster_pos]
50     return random.choice(available_positions)
51
52 def move(pos, action):
53     """
54     Given a position and an action, return the new position.
55     If the move is out-of-bounds, return the original position.
```

```

56     """
57     delta = actions_dict[action]
58     new_r = pos[0] + delta[0]
59     new_c = pos[1] + delta[1]
60     if 0 <= new_r < N and 0 <= new_c < N:
61         return (new_r, new_c)
62     else:
63         return pos
64
65 def step(state, agent_action, monster_action):
66     """
67     Execute one simultaneous step for agent and monster.
68     Returns next_state, reward, and a done flag.
69     """
70     agent_pos, monster_pos, apple_pos = state
71     new_agent_pos = move(agent_pos, agent_action)
72     new_monster_pos = move(monster_pos, monster_action)
73
74     # Check if agent and monster collide.
75     if new_agent_pos == new_monster_pos:
76         return (new_agent_pos, new_monster_pos, apple_pos), rewards_dict["caught_by_
monster"], True
77
78     # Check for apple collection.
79     reward = rewards_dict["empty"]
80     if new_agent_pos == apple_pos:
81         reward = rewards_dict["collect_apple"]
82         new_apple_pos = respawn_apple(new_agent_pos, new_monster_pos)
83     else:
84         new_apple_pos = apple_pos
85
86     next_state = (new_agent_pos, new_monster_pos, new_apple_pos)
87     return next_state, reward, False
88
89 def manhattan_distance(pos1, pos2):
90     """Return the Manhattan distance between two positions."""
91     return abs(pos1[0] - pos2[0]) + abs(pos1[1] - pos2[1])
92
93 # -----
94 # Feature Functions
95 # -----
96 def action_state_feature(state, action):
97     """
98     Computes a 4-dim feature vector x(s,a) given a state and an action.
99
100     Features:
101         f1 = 1/(distance from new agent position to apple + 1)
102         f2 = 1/(distance from new agent position to monster + 1)
103         f3 = 1 if the action moves the agent closer to the apple, else 0
104         f4 = 1 if the action moves the agent closer to the monster, else 0
105     """
106     agent_pos, monster_pos, apple_pos = state
107
108     # Current distances from agent position
109     current_dist_apple = manhattan_distance(agent_pos, apple_pos)
110     current_dist_monster = manhattan_distance(agent_pos, monster_pos)
111
112     # New agent position if action is taken

```

```

113 new_agent_pos = move(agent_pos, action)
114
115 # New distances from new agent position
116 new_dist_apple = manhattan_distance(new_agent_pos, apple_pos)
117 new_dist_monster = manhattan_distance(new_agent_pos, monster_pos)
118
119 f1 = 1 / (new_dist_apple + 1)
120 f2 = 1 / (new_dist_monster + 1)
121 f3 = 1 if new_dist_apple < current_dist_apple else 0
122 f4 = 1 if new_dist_monster < current_dist_monster else 0
123
124 return np.array([f1, f2, f3, f4])
125
126 def state_feature(state):
127     """
128     Returns a 4-dim state feature vector.
129     Here we use the differences between the agent and the apple/monster.
130     """
131     agent_pos, monster_pos, apple_pos = state
132     return np.array([apple_pos[0] - agent_pos[0],
133                     apple_pos[1] - agent_pos[1],
134                     monster_pos[0] - agent_pos[0],
135                     monster_pos[1] - agent_pos[1]])
136
137 # -----
138 # n-step SARSA (using action_state_feature)
139 # -----
140 def epsilon_greedy_theta(theta, state, epsilon):
141     """
142     Choose an action using the epsilon-greedy policy based on the linear approximator.
143     Here, theta is a vector and we use action_state_feature.
144     """
145     if np.random.rand() < epsilon:
146         return random.choice(actions)
147     q_vals = [np.dot(theta, action_state_feature(state, a)) for a in actions]
148     max_q = max(q_vals)
149     best_actions = [a for a, q in zip(actions, q_vals) if np.isclose(q, max_q)]
150     return random.choice(best_actions)
151
152 def n_step_sarsa(n, episodes=1000, alpha=0.001, gamma=0.95):
153     """
154     Implements semi-gradient n-step SARSA with a linear function approximator using
155     action_state_feature.
156     Theta is a 4-dim vector.
157
158     Returns:
159         theta: learned parameter vector.
160         episode_rewards: list of total returns per episode.
161     """
162     theta = np.zeros(4)
163     episode_rewards = []
164
165     for ep in range(episodes):
166         epsilon = get_epsilon(ep)
167         state = initial_positions() # (agent, monster, apple)
168         action = epsilon_greedy_theta(theta, state, epsilon)
169
170         states = [state]

```



```

170     actions_list = [action]
171     rewards = [0] # dummy for indexing
172
173     T_episode = float('inf')
174     t = 0
175     total_reward = 0
176     max_steps = T
177
178     while True:
179         if t < T_episode:
180             monster_action = random.choice(actions)
181             next_state, reward, done = step(states[t], actions_list[t], monster_
action)
182             rewards.append(reward)
183             total_reward += reward
184
185             if done:
186                 T_episode = t + 1
187             else:
188                 if t == max_steps - 1:
189                     T_episode = t + 1
190                 else:
191                     next_action = epsilon_greedy_theta(theta, next_state, epsilon)
192                     states.append(next_state)
193                     actions_list.append(next_action)
194
195             tau = t - n + 1
196             if tau >= 0:
197                 G = 0.0
198                 for i in range(tau + 1, min(tau + n, T_episode) + 1):
199                     G += (gamma ** (i - tau - 1)) * rewards[i]
200                 if tau + n < T_episode:
201                     feat = action_state_feature(states[tau+n], actions_list[tau+n])
202                     G += (gamma ** n) * np.dot(theta, feat)
203                 feat_tau = action_state_feature(states[tau], actions_list[tau])
204                 q_val = np.dot(theta, feat_tau)
205                 theta += alpha * (G - q_val) * feat_tau
206
207             if tau == T_episode - 1:
208                 break
209             t += 1
210
211             episode_rewards.append(total_reward)
212
213     return theta, episode_rewards
214
215 # -----
216 # REINFORCE and REINFORCE with Baseline (using state_feature for policy, but using
    action_state_feature for baseline)
217 # -----
218 # For the policy, we use the softmax over  $h(s,a) = \theta[:, a]^T \text{state\_feature}(s)$ .
219 def softmax_policy_reinforce(theta, state):
220     """
221     Computes the softmax policy using state features.
222     Theta is a (4 x |actions|) matrix.
223
224     Returns:
225     policy: dictionary mapping actions to probabilities.

```

```

226     phi: the state feature vector for state.
227     """
228     phi = state_feature(state)
229     scores = np.array([np.dot(theta[:, i], phi) for i, a in enumerate(actions)])
230     max_score = np.max(scores)
231     exp_scores = np.exp(scores - max_score)
232     probs = exp_scores / np.sum(exp_scores)
233     policy = {a: probs[i] for i, a in enumerate(actions)}
234     return policy, phi
235
236 def sample_action_reinforce(theta, state):
237     """Sample an action from the softmax policy (using state features)."""
238     policy, phi = softmax_policy_reinforce(theta, state)
239     chosen_action = np.random.choice(actions, p=[policy[a] for a in actions])
240     return chosen_action, policy, phi
241
242 def compute_returns(rewards, gamma):
243     """
244     Given a list of rewards (with a dummy at index 0), compute returns for each time
245     step.
246     """
247     G = 0
248     returns = []
249     for r in rewards[::-1]:
250         G = r + gamma * G
251         returns.insert(0, G)
252     return returns[1:] # skip dummy
253
254 def reinforce(epochs=1000, alpha=0.001, gamma=0.95):
255     """
256     Implements the REINFORCE algorithm using state features for the policy.
257     Theta is a (4 x |actions|) matrix.
258
259     Returns:
260         theta: learned policy parameters.
261         episode_rewards: list of total returns per episode.
262     """
263     theta = np.zeros((4, len(actions)))
264     episode_rewards = []
265
266     for ep in range(epochs):
267         state = initial_positions()
268         trajectory = [] # list of (state, action, reward)
269         t = 0
270         done = False
271
272         while not done and t < T:
273             action, policy, phi = sample_action_reinforce(theta, state)
274             monster_action = random.choice(actions)
275             next_state, reward, done = step(state, action, monster_action)
276             trajectory.append((state, action, reward))
277             state = next_state
278             t += 1
279
280         rewards_ep = [0] + [r for (_,_,r) in trajectory]
281         returns = compute_returns(rewards_ep, gamma)
282
283         for t_step, (state_t, action_t, _) in enumerate(trajectory):

```

```

283     policy_t, phi = softmax_policy_reinforce(theta, state_t)
284     # Update for each action in policy:
285     for i, a in enumerate(actions):
286         indicator = 1 if a == action_t else 0
287         grad = (indicator - policy_t[a]) * phi
288         theta[:, i] += alpha * returns[t_step] * grad
289     episode_rewards.append(sum([r for (_,_,r) in trajectory]))
290
291     return theta, episode_rewards
292
293 def reinforce_with_baseline(epochs=1000, alpha=0.001, beta=0.001, gamma=0.95):
294     """
295     Implements the REINFORCE algorithm with a state-value baseline.
296
297     The policy is parameterized using state features as before, but the baseline is
298     now approximated as
299      $v(s) = \max_a [w^T x(s,a)]$ ,
300     where  $x(s,a) = \text{action\_state\_feature}(s,a)$  and  $w$  is a 4-dim vector.
301
302     Returns:
303     theta: learned policy parameter matrix.
304     w: learned baseline parameter vector.
305     episode_rewards: list of total returns per episode.
306     """
307     theta = np.zeros((4, len(actions))) # policy parameters (for softmax using state
308     features)
309     w = np.zeros(4) # baseline parameters (using action_state_feature)
310     episode_rewards = []
311
312     for ep in range(epochs):
313         state = initial_positions()
314         trajectory = []
315         t = 0
316         done = False
317
318         while not done and t < T:
319             action, policy, phi = sample_action_reinforce(theta, state)
320             monster_action = random.choice(actions)
321             next_state, reward, done = step(state, action, monster_action)
322             trajectory.append((state, action, reward))
323             state = next_state
324             t += 1
325
326         rewards_ep = [0] + [r for (_,_,r) in trajectory]
327         returns = compute_returns(rewards_ep, gamma)
328
329         for t_step, (state_t, action_t, _) in enumerate(trajectory):
330             # Compute baseline  $v(s) = \max_a [w^T x(s,a)]$  using action_state_feature
331             q_values = [np.dot(w, action_state_feature(state_t, a)) for a in actions]
332             v_s = max(q_values)
333             # Determine the maximizing action a* (using argmax)
334             a_star = actions[np.argmax(q_values)]
335             phi_baseline = action_state_feature(state_t, a_star)
336             A_t = returns[t_step] - v_s # advantage
337
338             # Policy update (remains as before, using state features)
339             policy_t, phi_policy = softmax_policy_reinforce(theta, state_t)
340             for i, a in enumerate(actions):

```

```

339         indicator = 1 if a == action_t else 0
340         grad = (indicator - policy_t[a]) * phi_policy
341         theta[:, i] += alpha * A_t * grad
342
343         # Baseline update using the feature vector corresponding to a*
344         w += beta * A_t * phi_baseline
345
346         episode_rewards.append(sum([r for (_,_,r) in trajectory]))
347
348     return theta, w, episode_rewards
349
350
351 # -----
352 # Softmax Policy for Actor-Critic (using action_state_feature)
353 # -----
354 def softmax_policy_actor_critic(theta, state):
355     """
356     Computes the softmax policy for a given state using action_state_feature.
357
358     Parameters:
359         theta: actor parameter vector (4-dim)
360         state: current state
361
362     Returns:
363         policy: a dictionary mapping each action to its probability.
364     """
365     scores = []
366     for a in actions:
367         # The logit for action a
368         score = np.dot(theta, action_state_feature(state, a))
369         scores.append(score)
370     # Numerical stability: subtract max score
371     max_score = np.max(scores)
372     exp_scores = [np.exp(s - max_score) for s in scores]
373     sum_exp = np.sum(exp_scores)
374     probs = [s/sum_exp for s in exp_scores]
375     policy = {a: probs[i] for i, a in enumerate(actions)}
376     return policy
377
378 # -----
379 # One-Step Actor-Critic Algorithm
380 # -----
381 def one_step_actor_critic(episodes=15000, alpha=0.01, beta=0.01, gamma=0.95):
382     """
383     Implements the one-step actor-critic algorithm.
384
385     The policy is defined as:
386     
$$p(a|s; \theta) = \frac{\exp(\theta^T x(s,a))}{\sum_b \exp(\theta^T x(s,b))}$$

387     with  $x(s,a) = \text{action\_state\_feature}(s,a)$ .
388
389     The critic estimates the state value as:
390     
$$V(s) = \max_a [w^T x(s,a)]$$

391     and is updated using the feature vector corresponding to the maximizing action.
392
393     Updates:
394     Actor:  $\theta = \theta + \alpha [x(s,a) - \sum_b p(b|s; \theta) x(s,b)]$ 
395     Critic:  $w = w + \gamma x(s,a^*), \text{ where } a^* = \text{argmax}_a [w^T x(s,a)]$ 

```

```

397
398 Returns:
399     theta: learned actor parameter vector (4-dim)
400     w: learned critic parameter vector (4-dim)
401     episode_rewards: list of total return per episode.
402 """
403 # Initialize actor and critic parameters as zero vectors
404 theta = np.zeros(4)
405 w = np.zeros(4)
406 episode_rewards = []
407
408 for ep in range(episodes):
409     state = initial_positions()
410     done = False
411     total_reward = 0
412     t = 0
413     while not done and t < T:
414         # Select action using the softmax policy
415         policy = softmax_policy_actor_critic(theta, state)
416         a = np.random.choice(actions, p=[policy[a] for a in actions])
417
418         # Execute action; monster acts randomly
419         monster_action = random.choice(actions)
420         next_state, reward, done = step(state, a, monster_action)
421         total_reward += reward
422
423         # Critic: compute  $V(s) = \max_a [w^T x(s,a)]$ 
424         q_values_state = [np.dot(w, action_state_feature(state, b)) for b in
actions]
425         V_s = max(q_values_state)
426         # For terminal state, set  $V(s') = 0$ 
427         if done:
428             V_next = 0.0
429         else:
430             q_values_next = [np.dot(w, action_state_feature(next_state, b)) for b
in actions]
431             V_next = max(q_values_next)
432
433         # TD error
434         delta = reward + gamma * V_next - V_s
435
436         # Actor update:
437         # Compute gradient:  $\log(a|s) = x(s,a) - \sum_b (b|s) x(s,b)$ 
438         grad_log = action_state_feature(state, a)
439         for b in actions:
440             grad_log -= policy[b] * action_state_feature(state, b)
441         theta = theta + alpha * delta * grad_log
442
443         # Critic update:
444         # Find the action that maximizes  $w^T x(s,a)$ 
445         a_star = actions[np.argmax(q_values_state)]
446         phi_star = action_state_feature(state, a_star)
447         w = w + beta * delta * phi_star
448
449         state = next_state
450         t += 1
451
452     episode_rewards.append(total_reward)

```

```

453         if (ep+1) % 1000 == 0:
454             print("Episode {}: Total Reward = {}".format(ep+1, total_reward))
455
456         return theta, w, episode_rewards
457
458     %%
459     # -----
460     # Utility: Rolling Average Function
461     # -----
462     def rolling_average(data, window_size):
463         return np.convolve(data, np.ones(window_size)/window_size, mode='valid')
464
465     # -----
466     # Run All Methods
467     # -----
468     # n-step SARSA for n=1,2,3 (using action_state_feature)
469     theta_1, rewards_1 = n_step_sarsa(n=1, episodes=episodes, alpha=alpha, gamma=gamma)
470     theta_2, rewards_2 = n_step_sarsa(n=2, episodes=episodes, alpha=alpha, gamma=gamma)
471     theta_3, rewards_3 = n_step_sarsa(n=3, episodes=episodes, alpha=alpha, gamma=gamma)
472
473     %%
474     # REINFORCE and REINFORCE with Baseline (policy uses state_feature; baseline uses
475     #   action_state_feature)
476     theta_reinf, rewards_reinf = reinforce(episodes=episodes, alpha=alpha, gamma=gamma)
477     theta_rb, w_rb, rewards_rb = reinforce_with_baseline(episodes=episodes, alpha=0.0,
478     beta=0.05, gamma=gamma)
479
480     %%
481     # actor-critic
482     theta_ac, w_ac, rewards_ac = one_step_actor_critic(episodes=episodes, alpha=alpha,
483     beta=alpha, gamma=gamma)
484
485     # -----
486     # Compute rolling averages (using a window size appropriate for episodes)
487     window = 500
488     roll_rewards_1 = rolling_average(rewards_1, window)
489     roll_rewards_2 = rolling_average(rewards_2, window)
490     roll_rewards_3 = rolling_average(rewards_3, window)
491     roll_reinf = rolling_average(rewards_reinf, window)
492     roll_rb = rolling_average(rewards_rb, window)
493     roll_rewards_ac = rolling_average(rewards_ac, window)
494
495     # -----
496     # Plot All Methods Together
497     # -----
498     plt.figure(figsize=(10,6))
499     plt.plot(roll_rewards_1, label='n-step SARSA (n=1)')
500     plt.plot(roll_rewards_2, label='n-step SARSA (n=2)')
501     plt.plot(roll_rewards_3, label='n-step SARSA (n=3)')
502     plt.plot(roll_reinf, label='REINFORCE')
503     plt.plot(roll_rb, label='REINFORCE with Baseline')
504     plt.plot(roll_rewards_ac, label='One-Step Actor-Critic')
505     plt.xlabel('Episode')
506     plt.ylabel('Rolling Average Reward per (last {} episodes)'.format(window))
507     plt.title('Comparison of RL Methods on Gridworld')
508     plt.legend()
509     plt.show()

```



## References

Sutton, Richard S (2018). “Reinforcement learning: An introduction”. In: *A Bradford Book*.