# Project 2: TD methods and n-step bootstrapping

August Jonasson

February 25, 2025

## Part 1: Expected SARSA and importance sampling ratio

Suppose environment consists of four states $\{s_0, s_1, s_2, s_3\}$, with an action space $\{+, -\}$ and deterministic dynamics $p(s_{i+1}|s_i, +) = p(s_{i-1}|s_i, -) = 1$, with exceptions $p(s_0|s_3, +) = p(s_3|s_0, -) = 1$.

### Task 1

Given policy $\pi(+) = \pi(-) = 0.5$, $\forall s$, and having observed the transition

$$(s_0, +) \rightarrow (s_1, r = 2),$$

we want to compute the one-step *expected SARSA* update of the action-value estimate $Q(s_0, +)$. In order to do this, we assume initial estimates $Q(s, a) = 1$, $\forall s, a$, learning rate $\alpha = 0.2$, and discount factor $\gamma = 0.9$. Now, using equation (6.9) from the course book (Sutton 2018, p. 133)

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

with the given policy, parameters, values and observations, we get the one-step update

$$Q(s_0, +) \leftarrow 1 + 0.2 \left[ 2 + 0.9(0.5 + 0.5) - 1 \right] = 1.38,$$

and we are done.

### Task 2

Suppose an agent is following a behavior policy $b(+) = 0.7$ and $b(-) = 0.3$, with the target policy equal to that of the previous task. Having observed the trajectory

$$(s_0, +) \rightarrow (s_1, +) \rightarrow s_3,$$

we want to compute the *importance sampling ratio* $\rho_{0:1}$. In order to do this, we simply use equation (5.3) from the course book (p. 104), which defines

$$\rho_{t:T-1} := \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}.$$

Plugging in our policies for the two time steps, we get

$$\rho_{0:1} = \left(\frac{0.5}{0.7}\right)^2 \approx 0.51,$$

which can be interpreted as: "the observed trajectory is roughly half as likely to occur under the target policy than under the behavior policy", and we are done.

## Part 2: TD methods applied to gridworld with a monster

In this task, an agent acts in an $N^2$ gridworld with a randomly moving monster and a randomly spawned, stationary apple. The goal of the agent is to collect the apple as many times as possible before the episode is over, without getting caught by the monster.

Each state is represented by three tuples (agent position, monster position and apple position) and the entire space is given by

$$\mathcal{S} = \left\{(x_p, y_p), (x_m, y_m), (x_a, y_a) \,\middle|\, x_p, y_p, x_m, y_m, x_a, y_a \in \{0, \dots N-1\}\right\},$$

i.e. $|\mathcal{S}| = N^6$. The action space is

$$\mathcal{A} = \{\text{left}, \text{up}, \text{right}, \text{down}\},$$

which yields that the number of state-action pairs is $4 \cdot N^6$.

Some dynamics of the system are:

- the agent moves deterministically according to chosen action; if a wall is hit, the agent instead remains in place,

- the monster moves uniformly random in all four directions; a wall hit is the same as for the agent,

- the apple, if collected by agent, randomly respawns in a new empty cell,

- an episode ends if the monster catches the agent, or after $T$ time steps.

Given a state-action pair, the rewards are

$$R(s, a) = \begin{cases} + & 1 & \text{if the agent collects an apple,} \\ - & 1 & \text{if the agent is caught by the monster,} \\ & 0 & \text{otherwise.} \end{cases}$$

For the following tasks, we use hyperparameters:

- side length $N = 5$,

- eipsode length $T = 30$,

- discont factor $\gamma = 0.9$,

- learning rate $\alpha = 0.1$ and

- all state-action value estimates $Q(s, a)$ initialized at zero.

## Task 1: write code

In this first task we implement the gridworld system as defined above, as well as the following TD methods, using an $\epsilon$-greedy policy, for action-value estimation: SARSA, off-policy Q-learning, Double Q-learning, and n-step SARSA.

See Appendix 1 for the Python code.

## Task 2: plot the learning curves

In this task we are to answer which method converges the fastest and we also want to decide which $n$ is optimal, in the $n$-step SARSA method.
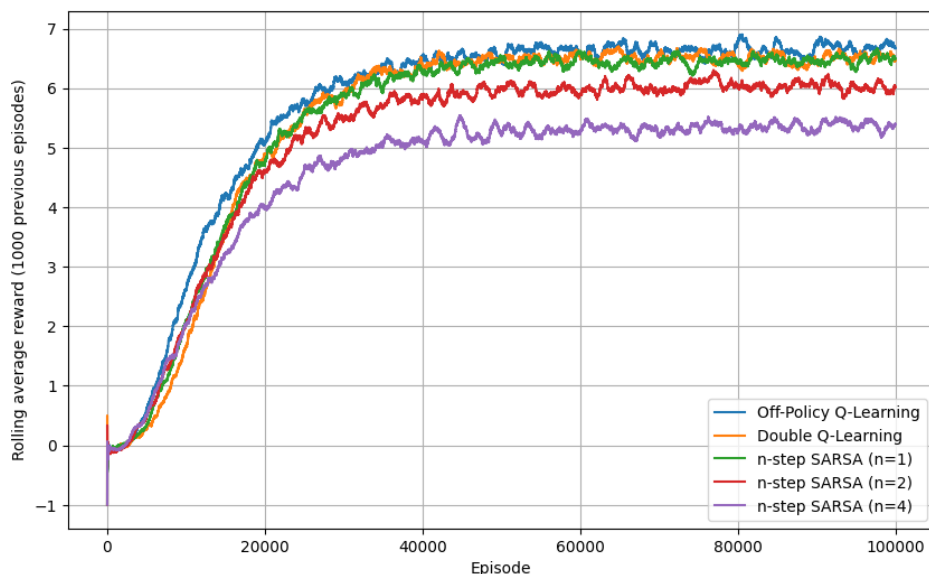


Figure 1: Moving average reward per episode (each average computed with previous 1000 episodes) for the different TD methods using $\epsilon$-greedy policy with exponentially decaying $\epsilon$ as a function of the current number of episodes (such that $\epsilon_1 = 1.0$ and $\epsilon_{10^5} \approx 0.01$).

From Figure 1 we see that all methods have similar convergence rates. For all five of them, convergence seems to happen after somewhere between 40 and 60 thousand episodes.

Considering the $n$-step SARSA methods, as we already mentioned, they converge at similar rates. Thus, if we are to answer which one is optimal we look to the one that produces the highest average reward, which is $n = 1$, and we are done.

3

## Task 3: discussion

In this last task, we discuss the following:

1. why all three $n$-step SARSA methods have similar rates of convergence,

2. why the apparent differences between regular Q-learning and Double Q-learning,

3. how convergence would be affected if the monster instead moves deterministically towards the agent.

**1)**

There is a case to be made for the bias-variance trade-off when it comes to these $n$-steps methods. A lower $n$ would correspond to a higher bias but low variance wheras a higher $n$ would have the reversed effects. Why we see so little variation in terms of rate of convergence among the three different $n$ could be that the trade-offs are pretty linear. Maybe it has something to do with the randomness of how the monster moves.

**2)**

The regular Q-learning method potentially suffers from the maximization bias (p. 134 in the course book), which means that we expect it to overestimate its action values early on in the training. The Double Q-learning method is the remedy to this bias and as such, it produces more reliable action values throughout the whole training. As we, in Figure 1, are looking at the actual rewards that result from the potentially biased policy and not the action-values themselves, the bias isn't apparent. We do, however, see a more rapid increase in reward early on from the Q-learning method, before the methods seem to converge at very similar values. This indicates that in this case, the maximization bias might actually be beneficial to early rewards.

**3)**

Convergence should be quicker as there are less paths to consider. Instead of $\approx 4$ monster moves per agent action, with a deterministically moving monster the number of paths should at least be halved. Thus leading to a smaller state-action space, thus leading to faster exploring of the entire space, thus leading to faster convergence.

# Appendix 1: Python code

```python
# %%
import numpy as np
import random
from matplotlib import pyplot as plt
import pandas as pd

# Gridworld parameters
N = 5                    # grid size (N x N)
T = 30                   # maximum time steps per episode
alpha = 0.1              # learning rate
gamma = 0.9              # discount factor
episodes = 10**5         # number of episodes per method

# Create grid coordinates
grid = [(r, c) for r in range(N) for c in range(N)]

# Define actions and their effects (using (row, col) convention)
actions = ["up", "right", "down", "left"]
actions_dict = {
    "up": (-1, 0),
    "right": (0, 1),
    "down": (1, 0),
    "left": (0, -1)
}

# Define rewards
rewards_dict = {
    "collect_apple": 1,
    "caught_by_monster": -1,
    "empty": 0
}

# Parameters for epsilon decay
epsilon_0 = 1.0   # Initial exploration rate
epsilon_min = 0.01 # Minimum exploration
decay_rate = 0.0001 # Controls speed of decay

def get_epsilon(episode):
    """ Exponentially decaying epsilon """
    return epsilon_min + (epsilon_0 - epsilon_min) * np.exp(-decay_rate * episode)

def initial_positions():
    """Return distinct starting positions for agent, monster, and apple."""
    return tuple(random.sample(grid, 3))

def respawn_apple(agent_pos, monster_pos):
    """Return a new apple position that is not occupied by agent or monster."""
    available_positions = [pos for pos in grid if pos != agent_pos and pos != monster_
    pos]
    return random.choice(available_positions)

def move(pos, action):
    """
    Given a position and an action, return the new position.
    If the move is out-of-bounds, return the original position.
    """
```

```python
56      delta = actions_dict[action]
57      new_r = pos[0] + delta[0]
58      new_c = pos[1] + delta[1]
59      if 0 <= new_r < N and 0 <= new_c < N:
60          return (new_r, new_c)
61      else:
62          return pos
63
64  def step(state, agent_action, monster_action):
65      """
66      Execute one simultaneous step for agent and monster.
67      Returns next_state, reward, and a done flag.
68      """
69      agent_pos, monster_pos, apple_pos = state
70      new_agent_pos = move(agent_pos, agent_action)
71      new_monster_pos = move(monster_pos, monster_action)
72
73      # Check if agent and monster collide.
74      if new_agent_pos == new_monster_pos:
75          return (new_agent_pos, new_monster_pos, apple_pos), rewards_dict["caught_by_
    monster"], True
76
77      # Check for apple collection.
78      reward = rewards_dict["empty"]
79      if new_agent_pos == apple_pos:
80          reward = rewards_dict["collect_apple"]
81          new_apple_pos = respawn_apple(new_agent_pos, new_monster_pos)
82      else:
83          new_apple_pos = apple_pos
84
85      next_state = (new_agent_pos, new_monster_pos, new_apple_pos)
86      return next_state, reward, False
87
88  def epsilon_greedy(Q, state, epsilon):
89      """
90      Return an action chosen by the epsilon-greedy policy based on Q.
91      """
92      if state not in Q or random.random() < epsilon:
93          return random.choice(actions)
94      max_val = max(Q[state].values())
95      best_actions = [a for a, v in Q[state].items() if v == max_val]
96      return random.choice(best_actions)
97
98  def epsilon_greedy_double(Q1, Q2, state, epsilon):
99      """
100     For double Q-learning: choose an action using the sum of Q1 and Q2 values.
101     """
102     if state not in Q1 or state not in Q2 or random.random() < epsilon:
103         return random.choice(actions)
104     combined = {a: Q1[state][a] + Q2[state][a] for a in actions}
105     max_val = max(combined.values())
106     best_actions = [a for a, v in combined.items() if v == max_val]
107     return random.choice(best_actions)
108
109 # %%
110 # --------------------------
111 # Off-Policy Q-Learning
112 # --------------------------
```

```python
113 Q_learning = {}
114 Q_learning_rewards = []
115
116 print("Starting Off-Policy Q-Learning...")
117 for episode in range(episodes):
118     state = initial_positions()  # (agent_pos, monster_pos, apple_pos)
119     if state not in Q_learning:
120         Q_learning[state] = {a: 0.0 for a in actions}
121
122     epsilon = get_epsilon(episode)
123     total_reward = 0
124     t = 0
125     done = False
126     while t < T and not done:
127         # Select action using Q_learning's own epsilon-greedy.
128         agent_action = epsilon_greedy(Q_learning, state, epsilon)
129         monster_action = random.choice(actions)
130
131         next_state, reward, done = step(state, agent_action, monster_action)
132         total_reward += reward
133
134         if not done:
135             if next_state not in Q_learning:
136                 Q_learning[next_state] = {a: 0.0 for a in actions}
137             best_next = max(Q_learning[next_state].values())
138         else:
139             best_next = 0
140
141         Q_learning[state][agent_action] += alpha * (reward + gamma * best_next - Q_
    learning[state][agent_action])
142         state = next_state
143         t += 1
144
145     # save accumulated episode reward
146     Q_learning_rewards.append(total_reward)
147
148     if (episode + 1) % 100 == 0:
149         print(f"Off-Policy Q-Learning Episode {episode+1}: Total Reward = {total_
    reward}")
150
151 # %%
152 # --------------------------
153 # 1-Step SARSA
154 # --------------------------
155 Q_sarsa = {}
156 Q_sarsa_rewards = []
157
158 print("\nStarting 1-Step SARSA...")
159 for episode in range(episodes):
160     state = initial_positions()
161     if state not in Q_sarsa:
162         Q_sarsa[state] = {a: 0.0 for a in actions}
163     # Choose initial action using Q_sarsa.
164     agent_action = epsilon_greedy(Q_sarsa, state, epsilon)
165
166     epsilon = get_epsilon(episode)
167     total_reward = 0
168     t = 0
```

```
169    done = False
170    while t < T and not done:
171        if state not in Q_sarsa:
172            Q_sarsa[state] = {a: 0.0 for a in actions}
173        monster_action = random.choice(actions)
174        next_state, reward, done = step(state, agent_action, monster_action)
175        total_reward += reward
176
177        if not done:
178            if next_state not in Q_sarsa:
179                Q_sarsa[next_state] = {a: 0.0 for a in actions}
180            next_action = epsilon_greedy(Q_sarsa, next_state, epsilon)
181            Q_sarsa[state][agent_action] += alpha * (reward + gamma * Q_sarsa[next_
    state][next_action] - Q_sarsa[state][agent_action])
182        else:
183            Q_sarsa[state][agent_action] += alpha * (reward - Q_sarsa[state][agent_
    action])
184
185        state = next_state
186        if not done:
187            agent_action = next_action
188        t += 1
189
190    Q_sarsa_rewards.append(total_reward)
191    if (episode + 1) % 100 == 0:
192        print(f"SARSA Episode {episode+1}: Total Reward = {total_reward}")
193
194 # --------------------------
195 # Double Q-Learning
196 # --------------------------
197 # In double Q-learning, the behavior policy is derived from the sum of Q_double1 and Q
    _double2.
198 Q_double1 = {}
199 Q_double2 = {}
200 Q_double_rewards = []
201
202 print("\nStarting Double Q-Learning...")
203 for episode in range(episodes):
204    state = initial_positions()
205    for Q in (Q_double1, Q_double2):
206        if state not in Q:
207            Q[state] = {a: 0.0 for a in actions}
208
209    epsilon = get_epsilon(episode)
210    total_reward = 0
211    t = 0
212    done = False
213    # Choose initial action using the double Q behavior policy.
214    agent_action = epsilon_greedy_double(Q_double1, Q_double2, state, epsilon)
215
216    while t < T and not done:
217        for Q in (Q_double1, Q_double2):
218            if state not in Q:
219                Q[state] = {a: 0.0 for a in actions}
220        monster_action = random.choice(actions)
221        next_state, reward, done = step(state, agent_action, monster_action)
222        total_reward += reward
223
```

```
224          for Q in (Q_double1, Q_double2):
225              if next_state not in Q:
226                  Q[next_state] = {a: 0.0 for a in actions}
227
228          if not done:
229              next_action = epsilon_greedy_double(Q_double1, Q_double2, next_state,
     epsilon)
230              # Randomly update one of the two Q-tables.
231              if random.random() < 0.5:
232                  best_action = max(Q_double1[next_state], key=Q_double1[next_state].get
     )
233                  target = reward + gamma * Q_double2[next_state][best_action]
234                  Q_double1[state][agent_action] += alpha * (target - Q_double1[state][
     agent_action])
235              else:
236                  best_action = max(Q_double2[next_state], key=Q_double2[next_state].get
     )
237                  target = reward + gamma * Q_double1[next_state][best_action]
238                  Q_double2[state][agent_action] += alpha * (target - Q_double2[state][
     agent_action])
239          else:
240              # Terminal state update.
241              if random.random() < 0.5:
242                  Q_double1[state][agent_action] += alpha * (reward - Q_double1[state][
     agent_action])
243              else:
244                  Q_double2[state][agent_action] += alpha * (reward - Q_double2[state][
     agent_action])
245
246          state = next_state
247          if not done:
248              agent_action = next_action
249          t += 1
250
251      Q_double_rewards.append(total_reward)
252      if (episode + 1) % 100 == 0:
253          print(f"Double Q-Learning Episode {episode+1}: Total Reward = {total_reward}")
254
255 # At the end, Q_learning, Q_sarsa, and (Q_double1, Q_double2) hold the learned action
     values for each method.
256
257 # %%
258 # ---------------------------
259 # n-step SARSA Function
260 # ---------------------------
261 def n_step_sarsa(n, episodes):
262      """
263      Runs n-step SARSA on the gridworld for a given number of episodes.
264      No pre-initialization of all states is needed; states are added as encountered.
265
266      Parameters:
267        n        : number of steps for bootstrapping
268        episodes : number of episodes to run
269
270      Returns:
271        Q_n      : learned Q-value table (dictionary)
272        rewards_n: list of total rewards per episode
273      """
```

```python
274     Q_n = {}
275     rewards_n = []
276     for episode in range(episodes):
277         epsilon_val = get_epsilon(episode)
278         # Initialize starting state and Q-values for that state if needed.
279         state = initial_positions()
280         if state not in Q_n:
281             Q_n[state] = {a: 0.0 for a in actions}
282         # Choose initial action using the epsilon-greedy policy.
283         action = epsilon_greedy(Q_n, state, epsilon_val)
284
285         # Lists to store the trajectory:
286         states = [state]            # states[0] is the initial state
287         actions_list = [action]     # actions_list[0] is the initial action
288         rewards_list = [0]          # rewards_list[0] is a dummy reward
289
290         T_episode = 30
291         t = 0
292         total_reward = 0
293
294         while True:
295             if t < T_episode:
296                 monster_action = random.choice(actions)
297                 next_state, reward, done = step(state, action, monster_action)
298                 total_reward += reward
299                 rewards_list.append(reward)
300                 states.append(next_state)
301                 if done:
302                     T_episode = t + 1
303                 else:
304                     if next_state not in Q_n:
305                         Q_n[next_state] = {a: 0.0 for a in actions}
306                     next_action = epsilon_greedy(Q_n, next_state, epsilon_val)
307                     actions_list.append(next_action)
308             tau = t - n + 1
309             if tau >= 0:
310                 # Determine the upper index for the reward sum.
311                 limit = min(tau + n, T_episode)
312                 G = 0.0
313                 for i in range(tau + 1, limit + 1):
314                     G += (gamma ** (i - tau - 1)) * rewards_list[i]
315                 if tau + n < T_episode:
316                     G += (gamma ** n) * Q_n[states[tau + n]][actions_list[tau + n]]
317                 Q_n[states[tau]][actions_list[tau]] += alpha * (G - Q_n[states[tau]][
    actions_list[tau]])
318             t += 1
319             if tau == T_episode - 1:
320                 break
321             # Update state and action only if available.
322             if t < len(states):
323                 state = states[t]
324             if t < len(actions_list):
325                 action = actions_list[t]
326         rewards_n.append(total_reward)
327         if (episode + 1) % 1000 == 0:
328             print(f"n-step SARSA (n={n}) Episode {episode+1}: Total Reward = {total_
    reward}")
329     return Q_n, rewards_n
```

```python
# %% [code]
# --------------------------
# Run n-step SARSA for different n values
# --------------------------
# print("\nStarting n-step SARSA (n=2)...")
_, rewards_n2 = n_step_sarsa(2, episodes)

# print("\nStarting n-step SARSA (n=4)...")
_, rewards_n4 = n_step_sarsa(4, episodes)

# For comparison, n=1-step SARSA is equivalent to standard 1-step SARSA:
print("\nStarting n-step SARSA (n=1) [equivalent to 1-step SARSA]...")
_, rewards_n1 = n_step_sarsa(1, episodes)


# %%
# --------------------------
# Plotting the Learning Curves
# --------------------------
avg_len = 1000

# creating lists of average over 1000 last episodes
Q_sarsa_avg = pd.Series(Q_sarsa_rewards).rolling(window=avg_len, min_periods=1).mean().to_numpy()
Q_learning_avg = pd.Series(Q_learning_rewards).rolling(window=avg_len, min_periods=1).mean().to_numpy()
Q_dbl_avg = pd.Series(Q_double_rewards).rolling(window=avg_len, min_periods=1).mean().to_numpy()

# Compute rolling averages over the last 1000 episodes
n1_avg = pd.Series(rewards_n1).rolling(window=avg_len, min_periods=1).mean().to_numpy()
n2_avg = pd.Series(rewards_n2).rolling(window=avg_len, min_periods=1).mean().to_numpy()
n4_avg = pd.Series(rewards_n4).rolling(window=avg_len, min_periods=1).mean().to_numpy()

episodes_range = list(range(1, episodes + 1))
plt.figure(figsize=(10, 6))
plt.plot(episodes_range, Q_learning_avg, label="Off-Policy Q-Learning")
# plt.plot(episodes_range, Q_sarsa_avg, label="1-Step SARSA")
plt.plot(episodes_range, Q_dbl_avg, label="Double Q-Learning")
plt.plot(episodes_range, n1_avg, label="n-step SARSA (n=1)")
plt.plot(episodes_range, n2_avg, label="n-step SARSA (n=2)")
plt.plot(episodes_range, n4_avg, label="n-step SARSA (n=4)")
plt.xlabel("Episode")
plt.ylabel("Rolling average reward (1000 previous episodes)")
plt.legend()
plt.grid(True)
plt.show()
# %%
```

# References

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