

COMP3667: Reinforcement learning practical 3

Markov decision processes

robert.lieck@durham.ac.uk & christopher.g.willcocks@durham.ac.uk

This is the version *with* answers!

1 Overview

Welcome to the third reinforcement learning practical. In this practical, we will try to better understand the concepts introduced in the last lecture by implementing and solving some of them using basic Python/NumPy. We will:

- Define and sample from a simple Markov chain.
- Extend to a Markov reward process and estimate approximate and compute exact state values.
- Extend to a Markov decision process and explore the influence of the policy.

2 Markov Chain

- Using Python and NumPy, define a Markov chain with $n_S = 10$ states $\mathcal{S} = \{0, \dots, n_S - 1\}$ and a transition matrix \mathcal{P} such that there is a probability of 0.9 for transitioning up ($s \rightarrow s + 1$) and 0.1 for transitioning down ($s \rightarrow s - 1$), wrapping around for the highest/lowest state. Only for $s = 0$ (the lowest state) there should be a probability of 0.8 for staying there and 0.1 of transitioning to the next higher/lower state. Check the normalisation of \mathcal{P} .

Answer:

```
import numpy as np
import matplotlib.pyplot as plt

n_states = 10
idx = np.arange(n_states)
transition_matrix = np.zeros((n_states, n_states))
transition_matrix[idx, (idx - 1) % n_states] = 0.1
transition_matrix[idx, (idx + 1) % n_states] = 0.9
transition_matrix[0, 0] = 0.8
transition_matrix[0, 1] = 0.1
transition_matrix[0, -1] = 0.1
print(f"normalised: {np.all(transition_matrix.sum(axis=1)==1)}")

→ normalised: True
```

- *Technical side remark:* Why can we use exact comparison (`==`) here, while we should normally use something like `np.isclose(x, y)` to compare floating point numbers?

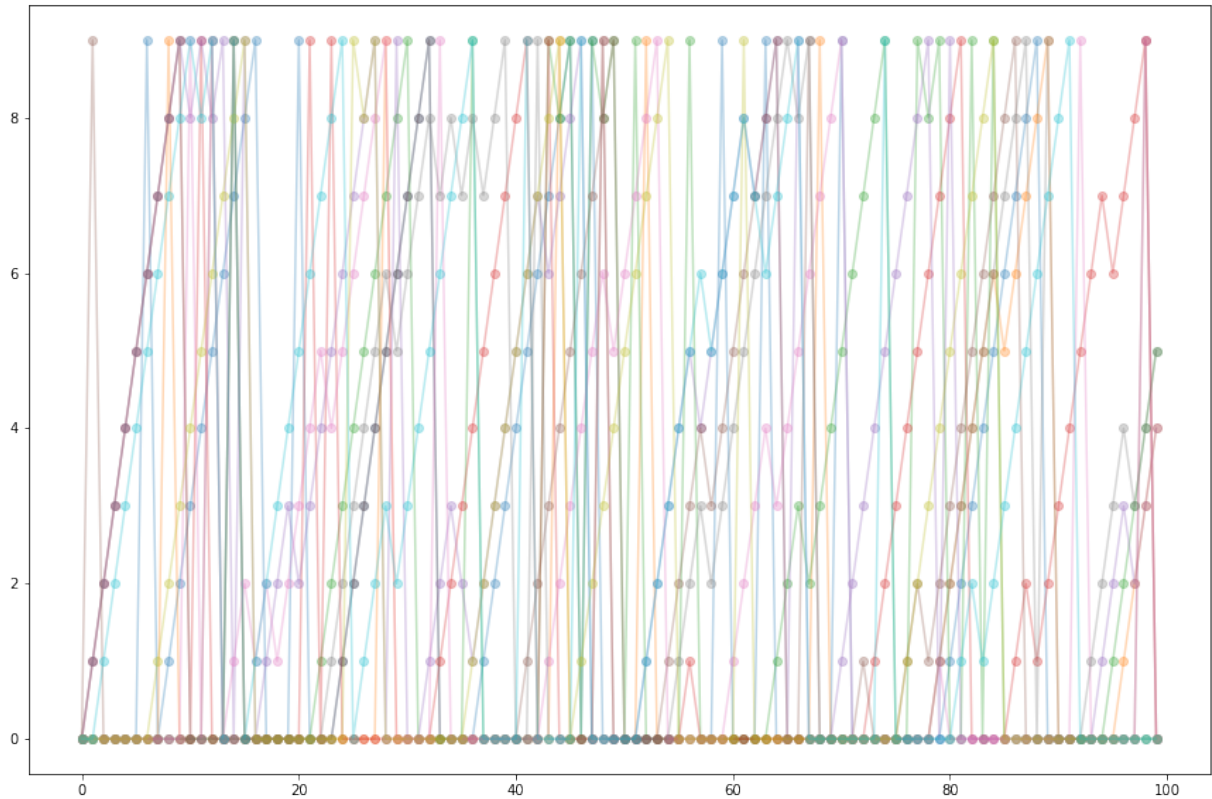
The initial values used here have an exact floating point representation and so do their sums, which we want to check.

- Write a function to sample n_E episodes of length l_E (always start in state $s = 0$). Generate episodes for $n_E = 10$ and $l_E = 100$ and plot them over time (i.e. each episode as a separate line with time on the horizontal axis and states on the vertical axis). What would you expect? Does the Markov chain behave as expected?

Answer:

```
def sample_episodes(n_episodes, length_episodes):
    episodes = np.zeros((n_episodes, length_episodes), dtype=int)
    for idx_episodes in range(n_episodes):
        for idx_time in range(1, length_episodes):
            current_state = episodes[idx_episodes, idx_time - 1]
            probs = transition_matrix[current_state]
            next_state = np.random.choice(n_states, p=probs)
            episodes[idx_episodes, idx_time] = next_state
    return episodes

episodes = sample_episodes(n_episodes=10, length_episodes=100)
fig, ax = plt.subplots(1, 1, figsize=(15, 10))
ax.plot(np.arange(episodes.shape[1]), episodes.T, '-o', alpha=0.3);
```



- Let $p(s_t)$ be the marginal distribution over the states at time t . Remember, a marginal distribution is when you do not know the other (potentially very informative/relevant) variables. So here, it is just a distribution over states, where we do not know what the previous state was (the previous state is marginalised out).
 - What is the marginal distribution over states at time $t + 1$?

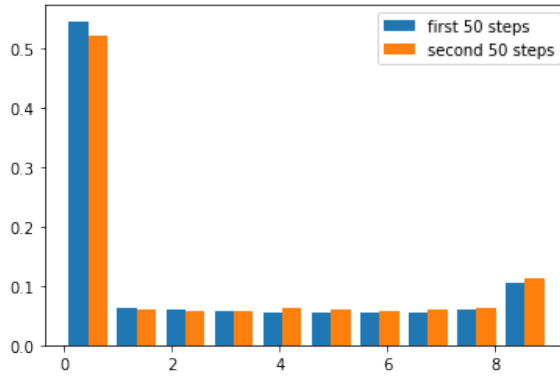
$$p(s_{t+1}) = \sum_{s_t} p(s_{t+1} | s_t) p(s_t) \quad (1)$$

- Sample another 100 episodes of length 100 and plot the marginal distribution over states
 - * for the first 50 time steps
 - * for the second 50 time steps.

You can do this using the `plt.hist` function after cutting all episodes appropriately and flattening them (two histogram plots can be done side by side by providing a tuple of data). Do you see a difference? Can you explain it?

As we start all episodes in state $s = 0$, the marginal distribution estimated from the first 50 steps has more probability mass on $s = 0$.

```
episodes = sample_episodes(n_episodes=100, length_episodes=100)
plt.hist((episodes[:, :50].flatten(), episodes[:, 50:].flatten()),
         density=True,
         label=["first 50 steps", "second 50 steps"])
plt.legend()
```



- **Advanced** [please, skip this question and come back to it later if you want, unless the solution is totally obvious to you and you love linear algebra]: What is the *stationary* distribution $\bar{p}(s)$ over states? This is defined by the following property: When you assume the stationary distribution as the marginal distribution at time t and compute the distribution at time $t + 1$, it is again the stationary distribution (it's stationary, it does not change). Mathematically speaking, it is an eigenvector of the transition operator \mathcal{P} with eigenvalue 1. Solve this by rewriting the condition as a matrix/vector equation and use `np.linalg.eig` to compute the eigenvalue and eigenvectors (note that `np.linalg.eig` computes right eigenvectors). Compare the result to your estimates of the marginal distribution above. Which one is more similar to the stationary distribution?

$$\bar{p}(s) \stackrel{!}{=} \sum_{s'} p(s | s') \bar{p}(s') \quad (2)$$

$$\Leftrightarrow \quad \bar{\mathbf{p}}_s \stackrel{!}{=} \bar{\mathbf{p}}_{s'} \mathbf{P}_{ss'} \quad (3)$$

$$\Leftrightarrow \quad \bar{\mathbf{p}}_s^\top \stackrel{!}{=} \mathbf{P}_{ss'}^\top \bar{\mathbf{p}}_{s'}^\top \quad (4)$$

```
eigvals, eigvecs = np.linalg.eig(transition_matrix.T)
stationary = None
for idx in range(n_states):
    if np.isclose(eigvals[idx], 1):
        stationary = eigvecs[:, idx].real
        stationary /= stationary.sum()
print(f"Is stationary: {np.all(np.isclose(stationary @ transition_matrix, stationary))}")
print(stationary)

--> Is stationary: True
--> [0.47368421 0.05263158 0.05263159 0.05263167 0.05263237 0.05263871
     0.05269575 0.05320915 0.05782976 0.0994152 ]
```

The estimate from the second 50 time steps is more similar to the stationary distribution. The time you need to wait for a Markov chain to produce reliable samples from the stationary distribution is called *mixing time*, which is very relevant in Markov chain Monte-Carlo methods.

3 Markov Reward Process

- Define a reward function that returns a reward of 1 when the chain enters state $s = 0$. Adapt your sampling function from above to also return rewards (along with the state at every time step in all episodes). Sample 1000 episodes of length 100 and plot the average reward for each time step. Does it converge? To what value? Compare to the probability of being in state 0 in the stationary distribution (0.47).

Answer:

```
def reward_function(s):
    return int(s == 0)

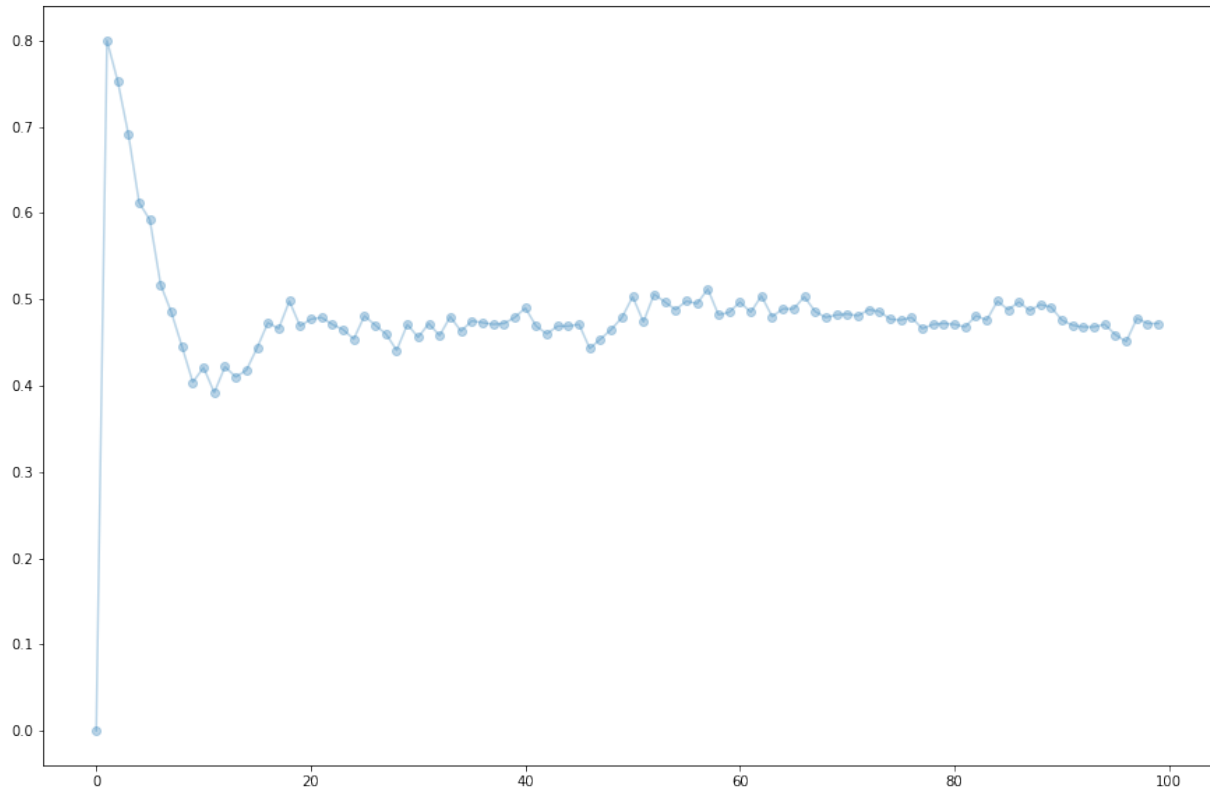
def sample_episodes(n_episodes, length_episodes):
    episodes = np.zeros((n_episodes, length_episodes, 2), dtype=int)
    for idx_episodes in range(n_episodes):
        for idx_time in range(1, length_episodes):
            current_state = episodes[idx_episodes, idx_time - 1, 0]
            probs = transition_matrix[current_state]
            next_state = np.random.choice(n_states, p=probs)
            episodes[idx_episodes, idx_time, 0] = next_state
```

```

        episodes[idx_episodes, idx_time, 1] = reward_function(next_state)
    return episodes

episodes = sample_episodes(n_episodes=1000, length_episodes=100)
fig, ax = plt.subplots(1, 1, figsize=(15, 10))
ax.plot(np.arange(episodes.shape[1]), episodes[:, :, 1].mean(axis=0), 'o', alpha=0.3);

```



- Assuming a discount of $\gamma = 0.5$, what is the value of being in states $s \in \{n_S - 1, 0, 1\}$? Estimate by sampling episodes from these states (adapt your sampling function to take an initial state).

Answer:

```

def sample_episodes(n_episodes, length_episodes, start_state):
    episodes = np.zeros((n_episodes, length_episodes, 2), dtype=int)
    episodes[:, 0, 0] = start_state
    for idx_episodes in range(n_episodes):
        for idx_time in range(1, length_episodes):
            current_state = episodes[idx_episodes, idx_time - 1, 0]
            probs = transition_matrix[current_state]
            next_state = np.random.choice(n_states, p=probs)
            episodes[idx_episodes, idx_time, 0] = next_state
            episodes[idx_episodes, idx_time, 1] = reward_function(next_state)
    return episodes

def estimate_state_values(*, discount, start_state_list, **kwargs):
    state_values = {}
    for start_state in start_state_list:
        episodes = sample_episodes(start_state=start_state, **kwargs)
        discounted_rewards = episodes[:, :, 1] * np.power(discount, np.arange(episodes.shape[1]))
        returns = discounted_rewards.sum(axis=1)
        value = returns.mean()
        state_values[start_state] = value
    print(f"Value for state {start_state}: {value}")
    return state_values

state_values = estimate_state_values(discount=0.5,
                                     start_state_list=[n_states - 1, 0, 1],
                                     n_episodes=100,
                                     length_episodes=10)

--> Value for state 0: 0.73376953125
--> Value for state 1: 0.10298828125
--> Value for state 9: 0.80435546875

```

- In the lecture we had the state value defined as

$$v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s'), \quad (5)$$

where \mathcal{R}_s is the reward for *being* in state s . Here we get a reward for *transitioning* to a state ($s = 0$). How would you need to adapt the equation for the state value to account for that?

$$v(s) = \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} [\mathcal{R}_{s'} + v(s')] \quad (6)$$

- Manually check the estimated state value for state $s = 0$. Is the Bellman quality approximately fulfilled?

$$v(s=0) = \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} [\mathcal{R}_{s'} + v(s')] \quad (7)$$

$$= 0.5 \left[\underbrace{0.1(0 + 0.80)}_{\text{to state } n_S-1} + \underbrace{0.8(1 + 0.73)}_{\text{stay in state 0}} + \underbrace{0.1(0 + 0.10)}_{\text{to state 1}} \right] \quad (8)$$

$$= 0.5 [0.08 + 1.384 + 0.01] = 0.5 \cdot 1.474 = \mathbf{0.737} \quad (9)$$

- The exact state values can be computed by rewriting the Bellman equation in matrix form and solving it (as in the lecture but again adapted because we get a reward not for *being in* but for *transitioning to a state*):

$$v = \gamma \mathcal{P}(\mathcal{R} + v) \quad (10)$$

$$\Leftrightarrow (\mathbf{I} - \gamma \mathcal{P})v = \gamma \mathcal{P}\mathcal{R} \quad (11)$$

Use `np.linalg.solve` to compute the exact state values.

Answer:

```
R = np.zeros(n_states)
R[0] = 1
np.linalg.solve(np.eye(n_states) - 0.5 * transition_matrix, 0.5 * transition_matrix @ R)
→ [0.74105284, 0.09069124, 0.00808578, 0.00789158, 0.01663844,
    0.03609746, 0.07836786, 0.17013997, 0.36938128, 0.80194284]
```

- The value for state $s = 0$ that you (a) estimated by sampling and (b) computed for checking the estimate both did not exactly match the true state value. Was the computed value (b) better or worse than the estimate from sampling (a)?

Most likely, your computed value (b) is closer to the true value than your initial estimate (a). In fact, when you use this to repeatedly update the values for *all* states in turn, this is guaranteed to converge to the true values (we will cover this method in the next lecture, it is called *value iteration*).

4 Markov Decision Process (MDP)

- Add two actions “up” and “down” to your Markov reward process to obtain an MDP. Define the transition matrix such that for action $a = 0$ (“up”) there is a 0.9 probability of transitioning to the next higher state and a 0.1 probability of transitioning to the next lower state (wrapping around in both directions). For action $a = 1$ (“down”) it should be the other way around.

Answer:

```
n_actions = 2 # up/down
idx = np.arange(n_states)
transition_matrix = np.zeros((n_states, n_states, n_actions))
transition_matrix[idx, (idx - 1) % n_states, 0] = 0.1
transition_matrix[idx, (idx + 1) % n_states, 0] = 0.9
transition_matrix[idx, (idx - 1) % n_states, 1] = 0.9
transition_matrix[idx, (idx + 1) % n_states, 1] = 0.1
print(f"normalised: {np.all(transition_matrix.sum(axis=1)==1)}")
→ normalised: True
```

- Extend your function for sampling episodes accordingly by passing the policy as an additional parameter and include the actions along with the state and reward.

Answer:

```
def sample_episodes(n_episodes, length_episodes, start_state, policy):
    episodes = np.zeros((n_episodes, length_episodes, 3), dtype=int)
    episodes[:, 0, 0] = start_state
    for idx_episodes in range(n_episodes):
        for idx_time in range(1, length_episodes):
            current_state = episodes[idx_episodes, idx_time - 1, 0]
            action = np.random.choice(n_actions, p=policy[current_state])
            probs = transition_matrix[current_state, :, action]
            next_state = np.random.choice(n_states, p=probs)
            episodes[idx_episodes, idx_time, 0] = next_state
            episodes[idx_episodes, idx_time, 1] = reward_function(next_state)
            episodes[idx_episodes, idx_time, 2] = action
    return episodes
```

- Define a uniform policy, which takes both actions with probability 0.5. Estimate the state value for states $s \in \{n_S - 1, 0, 1\}$ by sampling episodes under the uniform policy. Do the resulting values make sense? *Note:* If you defined your sampling function above in a generic way, you can reuse it unchanged.

Answer:

```
policy = np.zeros((n_states, n_actions))
policy[:, 0] = 0.5
policy[:, 1] = 0.5

estimate_state_values(discount=0.5,
                      start_state_list=[n_states - 1, 0, 1],
                      n_episodes=100,
                      length_episodes=10,
                      policy=policy);

--> Value for state 0: 0.1512109375
--> Value for state 1: 0.301015625
--> Value for state 9: 0.33361328125
```

- Change the policy to always go one step up and re-estimate the state values. Why is $v(s = n_S - 1)$ less than 0.5?

Answer:

```
policy = np.zeros((n_states, n_actions))
policy[:, 0] = 1
policy[:, 1] = 0

estimate_state_values(discount=0.5,
                      start_state_list=[n_states - 1, 0, 1],
                      n_episodes=100,
                      length_episodes=10,
                      policy=policy);

--> Value for state 9: 0.4971484375
--> Value for state 0: 0.050625
--> Value for state 1: 0.0397265625
```

- How would you change the policy to make it better? Try out and re-estimate, try to achieve the highest values for states $s \in \{n_S - 1, 0, 1\}$.

Answer:

```
policy = np.zeros((n_states, n_actions))
policy[:int(n_states / 2), 0] = 0
policy[:int(n_states / 2), 1] = 1
policy[int(n_states / 2):, 0] = 1
policy[int(n_states / 2):, 1] = 0

estimate_state_values(discount=0.5,
                      start_state_list=[n_states - 1, 0, 1],
                      n_episodes=100,
                      length_episodes=10,
                      policy=policy);

--> Value for state 9: 0.6121875
--> Value for state 0: 0.2923828125
--> Value for state 1: 0.610390625
```
