

Exercise 2: Markov Decision Processes (MDPs)

Please remember the following policies:

- Exercise due at **11:59 PM EST Feb 02, 2026.**
- Submissions should be made electronically on Canvas. Please ensure that your solutions for both the written and programming parts are present. You can upload multiple files in a single submission, or you can zip them into a single file. You can make as many submissions as you wish, but only the latest one will be considered.
- For **Written** questions, solutions may be handwritten or typeset. If you write your answers by hand and submit images/scans of them, please ensure legibility and order them correctly in a single PDF file.
- The PDF file should also include the figures from the **Plot** questions.
- For both **Plot** and **Code** questions, submit your source code in Jupyter Notebook (.ipynb file) along with reasonable comments of your implementation. Please make sure the code runs correctly.
- You are welcome to discuss these problems with other students in the class, but you must understand and write up the solution and code yourself. Also, you *must* list the names of all those (if any) with whom you discussed your answers at the top of your PDF solutions page.
- Each exercise may be handed in up to two days late (24-hour period), penalized by 10% per day late. Submissions later than two days will not be accepted.
- Contact the teaching staff if there are medical or other extenuating circumstances that we should be aware of.
- **Notations:** RL2e is short for the reinforcement learning book 2nd edition. x.x means the Exercise x.x in the book.
- Please note that question 4 is required for 5180 students and extra credit for 4180 students.

1. 1 point. *Formulating an MDP.*

Written: It is instructive to formally define an MDP at least once, in particular, the dynamics function. Consider the four-rooms domain from Ex0.

- What are the state and action spaces S, A ?
- Consider the dynamics function $p(s', r|s, a)$. Approximately how many rows in this conditional probability table have non-zero probability?
- 1 point. (Bonus, Optional) Code&Plot:** Provide codes to generate the full table of non-zero $p(s', r|s, a)$ values. Notes:
 - Please include both the code and the generated table of values.
 - Recall: The goal state is at $(10, 10)$. Any action from there teleports the agent back to $(0, 0)$.

2. 2 point. (RL2e 3.8, 3.9) *Discounted return.*

Written:

- Suppose $\gamma = 0.5$ and the following sequence of rewards is received:
 $R_1 = 1, R_2 = 5, R_3 = -2, R_4 = 4, R_5 = 4$
with $T = 5$. What are G_0, G_1, \dots, G_5 ?
- Suppose $\gamma = 0.9$ and the reward sequence is $R_1 = 3$ followed by an infinite sequence of ‘4’s.
What are G_1 and G_0 ?

3. 1 point. (RL2e 3.6, 3.7) *The RL objective.*

Written: Imagine that you are designing a robot to run a maze. You decide to give it a reward of +1 for escaping from the maze and a reward of zero at all other times. The task seems to break down naturally into episodes, where each episode corresponds to a single run through the maze and ends when agent escapes from the maze. The goal is to maximize expected total reward (Equation 3.7). There is no discounting, i.e. $\gamma = 1$. After running the learning agent for a while, you find that it is showing no improvement in escaping from the maze. What is going wrong? Have you effectively communicated to the agent what you want it to achieve?

1)

a) State(s) is the agents Position in the grid So its (x, y) coordinates in this example exo $(11 \times 11) = 121$ state
 but since walls can't be states $\# \text{walls} = 17$
 thus total states $= 121 - 17 = 104$, Action spec
 $= \{\begin{matrix} \text{up}, \text{down} \\ \text{left}, \text{right} \end{matrix}\}$

b) In exo there are 104 valid states
 And 4 actions. Thus $104 \times 4 = 416$

c) Code and plot is submitted as .ipyab file.

2)

$$y = 0.5, R_1 = 1, R_2 = 5, R_3 = -2$$

$$R_4 = 4, R_5 = 4$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} - \dots$$

$$G_5 = 0$$

$$G_3 = R_4 + 0.5 \times G_4$$

$$G_4 = R_5 = 4$$

$$= 4 + 0.5 \times 4 \\ = 6$$

$$G_2 = R_3 + 0.5 \times G_3 = -2 + 0.5 \times 6 \\ = 1$$

$$G_1 = R_2 + 0.5 + G_2 = 5 + 0.5 \times 1 \\ = 5.5$$

$$G_0 = R_1 + 0.5 \times G_1 = 1 + 0.5 \times 5.5 \\ = 3.25$$

(d)

$$G_1 = 4 + 0.9 \times 4 + 4 \times (0.9)^2, \dots$$

$$G_1 = 4 + \sum_{i=1}^{\infty} (0.9)^i = 4 \cdot \frac{1}{1-0.9} = 40$$

from geometric series

$$G_0 = R_1 + 0.9 \times G_1$$

$$= 3 + 0.9 \times 40 = 39$$

$$G_1 = 40 \quad G_0 = 39$$

3) with $\gamma = 0.99$ and reward = 1
only at escape, every policy
that escapes have same expected
(reward) which is 1 regardless
how efficient the policy is. Thus
agent has no learning signal
to prefer faster escape.

To solve this we could make
 $\gamma = 0.2$. to make
it more efficient.

4. **1 point.** (This question is required for 5180 students and extra credit for 4180 students) (RL2e 3.15, 3.16)
Modifying the reward function.

Written: In the gridworld example (Figure 3.2 in RL2e, see below), rewards are positive for goals, negative for running into the edge of the world, and zero the rest of the time. Are the signs of these rewards important, or only the intervals between them? Prove, using Equation 3.8, that adding a constant c to all the rewards adds a constant, v_c , to the values of all states, and thus does not affect the relative values of any states under any policies. What is v_c in terms of c and γ ?

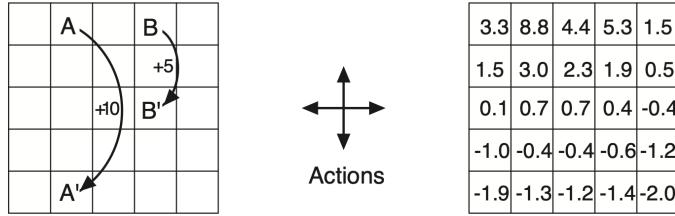


Figure 3.2: Gridworld example: exceptional reward dynamics (left) and state-value function for the equiprobable random policy (right).

5. **1 point.** (RL2e 3.14) *Bellman equation.*

Written: The Bellman equation holds for *all* policies, including optimal policies. Consider v_* and π_* shown in Figure 3.5 (middle, right respectively) (See below). Similar to the previous part, show numerically that the Bellman equation holds for the state positioned one bin to the right of the central state, valued at +16.0, with respect to its four neighboring states, for the optimal policy π_* shown in Figure 3.5 (right) (see below).

If the Bellman equation is unclear to you, you should practice this question again for other states.

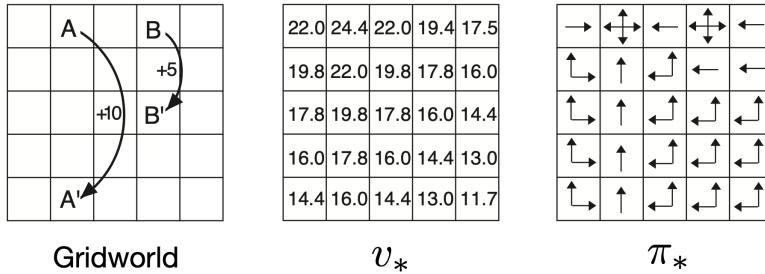


Figure 3.5: Optimal solutions to the gridworld example.

6. **2 points.** *Solving for the value function.*

Written: Another way to solve for the value function is to write out the Bellman equation for each state, view it as a system of linear equations, and then solve for the unknowns (the value of each state). Consider a particularly simple MDP, the 2-state recycling robot in Example 3.3.

- Expand the Bellman equation for the 2 states in the recycling robot, for an arbitrary policy $\pi(a|s)$, discount factor γ , and domain parameters $\alpha, \beta, r_{\text{search}}, r_{\text{wait}}$ as described in the example.
- You should now have two linear equations involving two unknowns, $v(\text{high})$ and $v(\text{low})$, as well as involving the policy $\pi(a|s)$, γ , and the domain parameters. Let $\alpha = 0.7, \beta = 0.5, \gamma = 0.9, r_{\text{search}} = 5, r_{\text{wait}} = 2$. Consider the policy $\pi(\text{search}|\text{high}) = 1, \pi(\text{wait}|\text{low}) = 0.4$, and $\pi(\text{recharge}|\text{low}) = 0.6$. Find the value function for this policy, i.e., solve the equations for the values of $v(\text{high})$ and $v(\text{low})$. Check that your solution satisfies the Bellman equation.
- 0.5 points. (Bonus, optional)** Suppose you can modify the policy in the `low` state, i.e., set $\pi(\text{wait}|\text{low}) = \theta$, and $\pi(\text{recharge}|\text{low}) = 1 - \theta$. What θ should you set it to, and what is the value function for that θ ?

$$4) \text{ eq. 3.9} \Rightarrow G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

After adding constant $\Rightarrow G_{t+k+1} = R_{t+k+1} + c$

$$= G_t + \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) \Rightarrow \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \sum_{k=0}^{\infty} \gamma^k c$$

First part $(\sum \gamma^k R_{t+k+1}) = G_t$

$$\Rightarrow G_t = G_t + (\sum \gamma^k c)$$

For $|y| < 1$ eq. 3.10

$$G_t = G_t + \frac{c}{1-y}$$



again

$$V_T(S) = E[G_t | S = s]$$

$$V_T(S) = V_P(S) + \frac{c}{1-y}$$

$$\text{thus } V_C = \frac{c}{1-y}$$

look

↓ Geometric series

So signs of rewards not important because adding constant to all rewards shifts values by same amount $(\frac{c}{1-y})$ (constant)

The difference between rewards are the important part.

5) Bellman eq.

$$\sum \pi(a|s) \sum p(s'|r/s,a) [r + \gamma v_{\pi}(s')]$$

Neighbours: up = 17.8, left = 17.8, right = 16.4
down = 14.4

Policy π_* = {left, up} of that state

thus? Split up, left equal probability

$$v_*(s) = 0.5 \times (0.9 \times 17.8) + 0.5 \times (0.9 \times 13.8)$$

$$\geq \frac{16.02 + 16.02}{2} = 16.02 \approx \underline{\underline{16}}$$

Assuming its deterministic and policy is as in the table (left, up)

b) Belman eq. $\sum \bar{\pi}(a|s) \sum p(s', r | s, a)$
 $[c + \gamma v_{\pi}(s')]$

a) From example 3.3

States = {High, Low}, Actions = $A(H) = \{Search, wait\}$

$A(L) = \{Search, wait, Recharge\}$

For High:

$$v(H) = \bar{\pi}(\text{Search}|H) \left(c_{\text{Search}} + \gamma (\alpha v(H) + (1-\alpha)v(L)) \right) \\ + \bar{\pi}(\text{wait}|H) \left(c_{\text{wait}} + \gamma v(H) \right)$$

for Low:

$$v(L) = \bar{\pi}(\text{Search}|L) \left(\beta ((\text{Search} + \gamma v(L)) + (1-\beta)(-3 + \gamma v(H))) \right) \\ + \bar{\pi}(\text{wait}|L) (c_{\text{wait}} + \gamma v(L)) + \bar{\pi}(\text{recharge}|L) (0 + \gamma v(H))$$

b) $\alpha = 0.7, \beta = 0.5, \gamma = 0.9, c_{\text{Search}} = 5, c_{\text{wait}} = 2$

Policy $\pi: \bar{\pi}(\text{Search}|H) = 1, \bar{\pi}(\text{wait}|L) = 0.2$

$$\bar{\pi}(\text{recharge}|L) = 0.6$$

High and always search

$$v(H) = 5 + 0.9 (0.7v(H) + 0.3v(L))$$

$$= 5 + 0.63v(H) + 0.27v(L) \rightarrow 0.37v(H) = 5 + 0.27v(L)$$

Low and wait / recharge mix!

$$V(L) = 0.1 (2 + 0.9 V(L)) + 0.6 (0 + 0.9 V(H))$$

$$V(L) = 0.8 + 0.36 V(H) + 0.54 V(H) \rightarrow 0.6 L + V(L) = 0.8 + 0.54 V(H)$$

$$V(L) = \frac{0.8 + 0.54 V(H)}{0.6} = 1.25 + 0.87 V(H)$$

$$V(H) \approx 33.54, V(L) \approx 32.92$$

$$c) \Theta = \text{P(wait | L)}, \text{P(recharge | L)} = 1 - \Theta$$

$$V(L) = \Theta (L + 0.9 V(L)) + (1 - \Theta) 0.9 V(H)$$

$$\text{if } \Theta = 0, V(L) = 0.9 V(H)$$

$$V(H) = 5 + 0.9 (0.7 V(H) + 0.3 V(L)) = 5 + 0.9 (0.7 V(H) + 0.27 V(H))$$

$$= 5 + 0.87 V(H)$$

$$0.127 V(H) \rightarrow V(H) = 39.37, V(L) \approx 35.43$$

