

# Linear Programming Assignment

Sandro Paradžik  
*University of Sarajevo*

December 23, 2024

## 1 Introduction

This document presents a solution to a Markov Decision Process (MDP) problem using linear programming (LP). The problem, as proposed by Sutton and Barto [1], involves optimizing the actions of a robot that collects soda cans while managing its battery levels efficiently.

## 2 Problem Statement

The robot operates in two battery states: **high** and **low**. Depending on the state, the robot can choose among several actions:

- **Search for cans:** Yields an expected reward of 2, but in the high state, it risks transitioning to the low state with probability  $1 - \alpha$ . In the low state, searching risks running out of battery, penalized by  $-3$  (after this battery is set to high state).
- **Wait for cans:** Provides an expected reward of 1 and keeps the battery state unchanged.
- **Charge the battery:** Available only in the low state, it transitions to the high state without a direct reward.

The objective is to maximize the cumulative discounted reward with a discount factor  $\gamma = 0.9$ .

## 3 Mathematical Formulation

Using the Bellman optimality equations [1], the problem is formulated as an LP [2]. Let  $v(h)$  and  $v(l)$  represent the value functions for the high and low battery states, respectively. The rewards are defined as  $r_{search} = 2$  and  $r_{wait} = 1$ . The constraints are derived as follows:

$$\begin{aligned} \text{High state (h): } & v(h) \geq r_{wait} + \gamma v(h), \\ & v(h) \geq r_{search} + \gamma (\alpha v(h) + (1 - \alpha)v(l)). \\ \text{Low state (l): } & v(l) \geq r_{wait} + \gamma v(l), \\ & v(l) \geq \gamma v(h), \\ & v(l) \geq \beta r_{search} - 3(1 - \beta) + \gamma ((1 - \beta)v(h) + \beta v(l)). \end{aligned}$$

The LP formulation is:

Minimize:  $v(h)$   
Subject to: the above constraints.

## 4 Python Implementation

The Python implementation uses the `cvxpy` library to solve the linear programming formulation of the recycling robot problem. This library simplifies the creation and solving of optimization problems with constraints. Below is the key idea behind the implementation:

- The `v_h` and `v_l` variables represent the value functions for the high and low battery states, respectively.
- The constraints are derived from the Bellman equations for each state, ensuring the policy satisfies the optimality conditions.
- The objective is to minimize `v_h`, which represents the value of the high battery state.
- The implementation includes logic to determine the optimal policy ( $\pi_*(h)$  and  $\pi_*(l)$ ) based on the calculated optimal value functions.

```
1  def recycling_robot(alpha, beta, r_s=2, r_w=1, gamma=0.9):
2      # Decision variables
3      v_h = cp.Variable(name="v_h") # Value for high state
4      v_l = cp.Variable(name="v_l") # Value for low state
5
6      # Objective
7      objective = cp.Minimize(v_h) # we can also use v_h + v_l
8
9      # Constraints (Bellman)
10     constraints = [
11         # high -> wait
12         v_h >= r_w + gamma*v_h,
13
14         # high -> search
15         v_h >= r_s + gamma*(alpha*v_h + (1 - alpha)*v_l),
16
17         # low -> wait
18         v_l >= r_w + gamma*v_l,
19
20         # low -> recharge
21         v_l >= gamma*v_h,
22
23         # low -> search
```

```

24         v_l >= beta*r_s - 3*(1 - beta) + gamma*((1 -
           beta)*v_h + beta*v_l)
25     ]
26
27     # Solve the problem using a linear programming solver
28     prob = cp.Problem(objective, constraints)
29     prob.solve(solver=cp.GLPK) # Use GLPK, an LP solver
30
31     # Convert v_h and v_l to float
32     v_h = float(v_h.value)
33     v_l = float(v_l.value)
34
35     # Calculate optimal policies
36     pi_h = -1
37     pi_l = -1
38
39     eps = 0.001
40
41     if abs(v_h - (r_w + gamma*v_h)) < eps:
42         pi_h = 1 # wait
43     elif abs(v_h - (r_s + gamma*(alpha*v_h + (1 -
           alpha)*v_l))) < eps:
44         pi_h = 2 # search
45
46     if abs(v_l - (r_w + gamma*v_l)) < eps:
47         pi_l = 1 # wait
48     elif abs(v_l - (gamma*v_h)) < eps:
49         pi_l = 0 # recharge
50     elif abs(v_l - (beta*r_s - 3*(1 - beta) + gamma*((1 -
           beta)*v_h + beta*v_l))) < eps:
51         pi_l = 2 # search
52
53     return {
54         "v_h": v_h,
55         "v_l": v_l,
56         "pi_h": pi_h,
57         "pi_l": pi_l
58     }

```

Listing 1: Solving the problem as an LP in Python.

## 5 Results and Discussion

We presented an approach to solving the recycling robot problem using Python. Figure 1 illustrates the optimal value functions and policies for  $\alpha \in (0, 1)$  and  $\beta \in (0, 1)$ , demonstrating how MDPs can be effectively solved using LP.

Another aspect of interest is the time complexity of this approach. While larger instances of MDPs are not typically solved using LP, the approach we used has strong theoretical guarantees. Specifically, LP can solve MDPs in polynomial time with respect to  $|S| \cdot |A|$ , where  $S$  represents the set of possible states and  $A$  the set of possible actions.

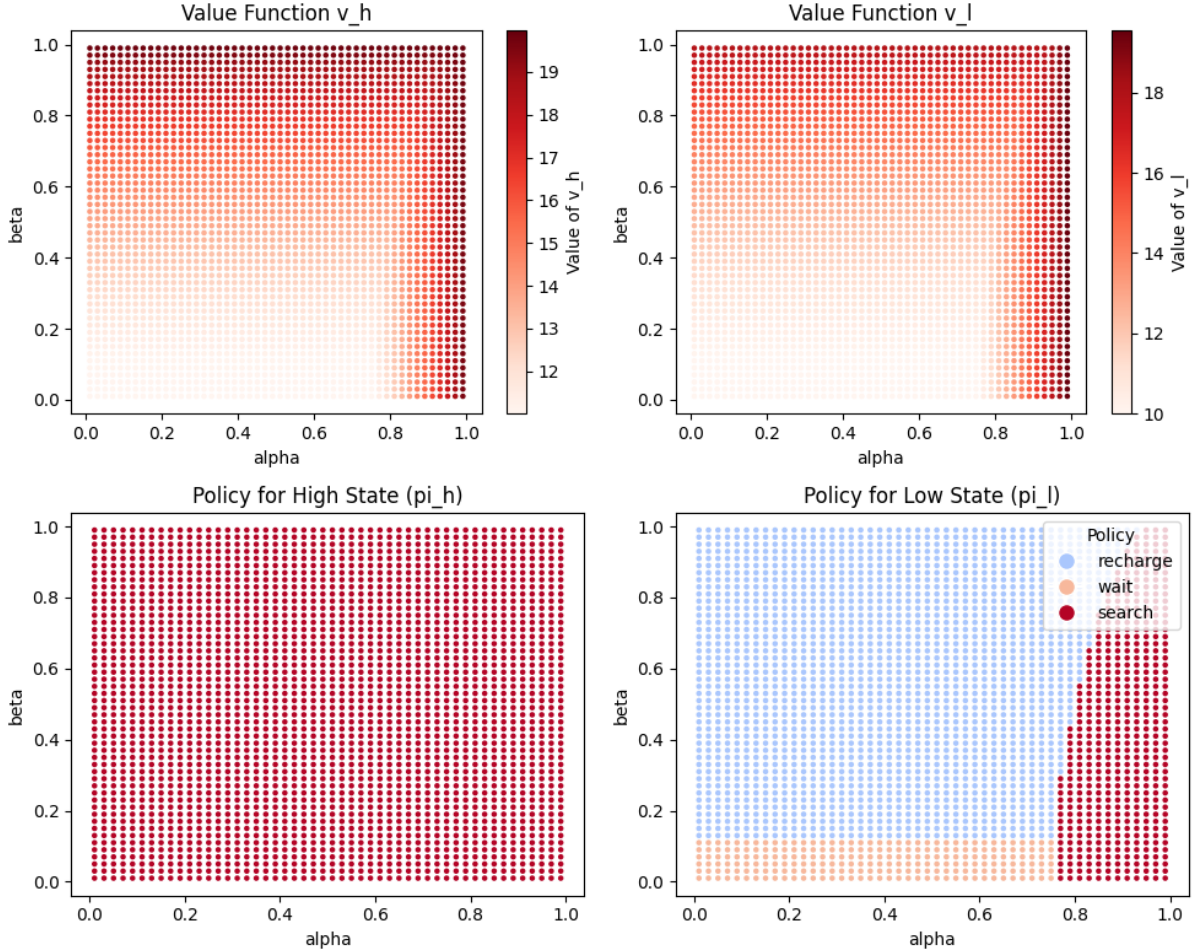


Figure 1: Results of the linear programming solution.

## References

- [1] Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction* (2nd ed.). MIT Press.
- [2] Helmert, M., & Röger, G. (2021). *Planning and Optimization: F2. Bellman Equation & Linear Programming*. Retrieved from [https://ai.dmi.unibas.ch/\\_files/teaching/hs21/po/slides/po-f02.pdf](https://ai.dmi.unibas.ch/_files/teaching/hs21/po/slides/po-f02.pdf)