

In this assignment, you solve problems for Markov Decision Processes (MDP) - Average. You need to formulate the mathematical model and solve it with Python and Gurobi.

1. Submission Instructions

Submit a PDF file describing a Markov Decision Processes (MDP) model and reporting the solution to the problem instance. Also, submit a program (a Python script or a Jupyter notebook) using Gurobi to solve the problem instance. In the MDP formulation, clearly define the stages, states, and actions and formulate the Bellman equations. In the Linear Programming (LP) formulation, clearly define decision variables and state the objective function and constraints. For the problem instance, report the values of the objective and solution.

2. Problem

A patient is suffering from a chronic disease that is non-lethal but incurable. Each stage not only impacts the patient's quality of life but also incurs increasing costs for symptom relief per cycle: \$100 for mild (M), \$300 for intermediate (I), and \$500 for severe (S). To manage this condition, two treatment options are available: basic medication (B) and advanced medication (A). Basic one costs \$100 per cycle, while the advanced costs \$400 but is more effective.

Treatment Effects:

- Mild Stage:**
Basic Medication: 80% chance of remaining mild, 20% chance of worsening to intermediate.
Advanced Medication: 90% chance of remaining mild, 10% chance of worsening to intermediate.
- Intermediate Stage:**
Basic Medication: 60% chance of staying intermediate, 40% chance of worsening to severe.
Advanced Medication: 70% chance of staying intermediate, 20% chance of improving to mild, 10% chance of worsening to severe.
- Severe Stage:**
Basic Medication: 100% chance of remaining severe.
Advanced Medication: 80% chance of staying severe, 15% chance of improving to intermediate, and 5% chance of improving to mild.

Question 1: draw the State Transition Diagram for MDP, formulate an Average Reward MDP as LP and then solve it using Python and Gurobi.

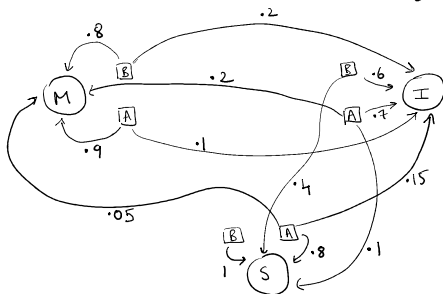
Question 2: formulate the Bellman equations for Discounted Reward MDP (discount factor is 0.9), perform value iteration with Python (100 iterations to converge) and report the optimal policy and value. Note that the solution to Problem 2 is NOT necessarily the same as Problem 1.

} SOLUTION BELOW AND ON GURUBI ↓

} SOLUTION BELOW AND ON PYTHON ↓

Problem 1 (Average Reward MDP)

State Transition Diagram : M = mild, I = intermediate, S = Severe
B = basic, A = advanced.



General Formulation

$$\text{Min } \sum_s \sum_x r_{sx} \pi_{sx}$$

$$\sum_s \sum_x \pi_{sx} = 1, \pi_{sx} \geq 0$$

$$\sum_x \pi_{sx} = \sum_x \sum_i \pi_{ix} P_{isx} \quad \forall s$$

→ transition state probabilities : see diagram

π_{sx} is the state and choice

$s \in \{m, i, s\}$: mild, intermediate, severe
 $x \in \{b, a\}$: basic, advanced

Decision variables : $(\pi_{sa}, \pi_{sb}, \pi_{ma}, \pi_{mb}, \pi_{ia}, \pi_{ib})$
to understand which states form optimal policy.

r_{sx} = fixed cost for symptom relief + treatment

$$r_{mb} = 200, r_{ib} = 400, r_{sb} = 600$$

... and

Linear cost policy

$$r_{mb} = 200, r_{ib} = 400, r_{sb} = 600$$

$$r_{ma} = 500, r_{ia} = 700, r_{sa} = 900$$

Objective

$$\min (r_{mb}\pi_{mb} + r_{ma}\pi_{ma} + r_{ib}\pi_{ib} + r_{ia}\pi_{ia} + r_{sb}\pi_{sb} + r_{sa}\pi_{sa}) \quad \text{minimize avg cost}$$

Constraints

$$\pi_{mb} + \pi_{ma} + \pi_{ib} + \pi_{ia} + \pi_{sb} + \pi_{sa} = 1 \quad \text{ensure state probabilities sum to 1, } \geq 0$$

$$\pi_{mb}, \pi_{ma}, \pi_{ib}, \pi_{ia}, \pi_{sb}, \pi_{sa} \geq 0$$

$$\begin{aligned} \pi_{mb} + \pi_{ma} &= 0.8\pi_{mb} + 0.9\pi_{ma} + 0.2\pi_{ia} + 0.05\pi_{sa} \\ \pi_{sb} + \pi_{sa} &= \pi_{sb} + 0.8\pi_{sa} + 0.1\pi_{ia} + 0.4\pi_{ib} \end{aligned} \quad \left. \begin{array}{l} \text{Prob of being in state} \\ \text{Prob of going to state} \end{array} \right\}$$

$$\pi_{ib} + \pi_{ia} = 0.6\pi_{ib} + 0.7\pi_{ia} + 0.15\pi_{sa} + 0.1\pi_{ma} + 0.2\pi_{mb} \quad \text{Not needed extra constraint.}$$

Problem Instance Solution (using Gurobi)

$$\begin{aligned} \pi_{mb} &= 0.428571 \\ \pi_{ma} &= 0 \\ \pi_{ib} &= 0 \\ \pi_{ia} &= 0.380952 \\ \pi_{sb} &= 0 \\ \pi_{sa} &= 0.190476 \end{aligned}$$

Decision variable values.

optimal policy must include states that take values

so $\pi_{mb}, \pi_{ia}, \pi_{sa}$

Thus, $x^* = (0, 1, 1)$, i.e. $x^*(m) = b$, $x^*(i) = a$, $x^*(s) = a$

optimal solution.

Objective minimized: \$523.809524 \sim \text{avg cost.}

optimal policy = $\{0, 1, 1\} = \{b, a, a\}$

i.e. basic treatment if mild condition,
advanced treatment if intermediate or severe condition.

NOTE: in code : $\pi[0] = \pi_{mb}, \pi[1] = \pi_{ma}$
 $\pi[2] = \pi_{ib}, \pi[3] = \pi_{ia}$
 $\pi[4] = \pi_{sb}, \pi[5] = \pi_{sa}$

the values are mapped accordingly.

Problem 2: Discounted Reward MDP

Stages: time period (need to make decision at each stage)

States: $s \in \{m, i, s\}$ (mild, intermediate, or severe condition)

Decision: $x \in \{b, a\}$ (basic or advanced treatment)

General Bellman \rightarrow 0.9 discount rate

$$V_s = \min_x \left\{ r_{sx} + \beta \sum_j p_{sj} V_j \right\} \quad \forall s$$

\rightarrow transition probabilities given

$$V_s = \min_x \left\{ r_{sx} + \beta \sum_j p_{sj} V_j \right\} \quad \forall s$$

↳ transition probabilities given

fixed symptom relief + treatment cost

$$r_{sx} : r_{MB} = 200, r_{MA} = 500, r_{IB} = 400, r_{IA} = 700, r_{SB} = 600, r_{SA} = 900$$

Bellman Equations

$$V_M = \min_{x \in \{b, a\}} \left\{ r_{MB} + 0.9(0.8V_M + 0.2V_I), r_{MA} + 0.9(0.9V_M + 0.1V_I) \right\}$$

$$V_I = \min_{x \in \{b, a\}} \left\{ r_{IB} + 0.9(0.6V_I + 0.4V_S), r_{IA} + 0.9(0.7V_I + 0.2V_M + 0.1V_S) \right\}$$

$$V_S = \min_{x \in \{b, a\}} \left\{ r_{SB} + 0.9(V_S), r_{SA} + 0.9(0.8V_S + 0.15V_I + 0.05V_M) \right\}$$

Problem Instance Solution (in code basic=0, advanced=1)

After 100 iterations of value iteration,

Initial guess: $V_M^{(0)}, V_I^{(0)}, V_S^{(0)} = 0$.] → Action for initial guess.

$[b, b, b]$ ↳ ie. $x^*(M) = b, x^*(I) = b, x^*(S) = b$

Converged values

$$V_M^{(100)} = \$4173.9982$$

$$V_I^{(100)} = \$5381.8633$$

$$V_S^{(100)} = \$5999.8406$$

converged values
after 100 iterations

Optimal Policy

$[b, a, b]$ → basic treatment for mild condition

$x^*(M) = b$, advanced treatment for intermediate condition

$x^*(I) = a$, basic treatment for severe condition.

$x^*(S) = b$

optimal treatment plan.