

Markov Decision Processes and Reinforcement Learning

Day 2 — Optimal Value Functions and Classical Computational Methods

Luiz Felipe Martins

Department of Mathematics and Statistics
Cleveland State University

January 27, 2026

Outline

- 1 Review of Previous Presentation
- 2 Optimal value functions
- 3 Tabular Solution Methods
 - Dynamic Programming

Definition of Markov Decision Process

- A trajectory of an MDP consists of three stochastic sequences:
 - S_0, S_1, S_2, \dots : the states visited by the environment.
 - A_1, A_2, A_3, \dots the actions chosen by the agent.
 - R_0, R_1, R_2, \dots the rewards awarded to the agent.
- Model is specified by two deterministic functions:

$$p(s' | s, a) \quad \text{and} \quad R(s, s', a).$$

- Model evolution:

$$\mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a) = p(s' | s, a)$$

$$R_{t+1} = r(S_t, S_{t+1}, A_t)$$

Rewards, returns and state value functions

- The **reward** received at time $t + 1$ (after action A_t is selected) is denoted by R_{t+1} .
- The **total return** (or simply return) for a run is:

$$\sum_{t=0}^T \gamma^t R_{t+1}$$

- The expected return of a run is:

$$\mathbb{E} \left[\sum_{t=0}^T \gamma^t R_{t+1} \right]$$

Policies and value functions

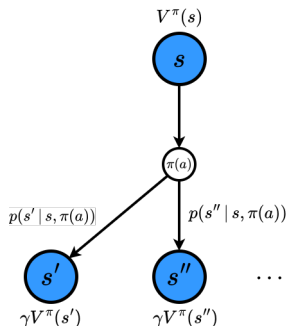
- A **deterministic policy** associates to each state s and action $a = \pi(s)$, where a is in the set of admissible actions $\mathcal{A}(s)$
- The state value function for policy π is:

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^T \gamma^t R_{t+1} \mid S_0 = s \right]$$

- We assume:
 - For episodic tasks, $P(T < \infty) = 1$ and $0 < \gamma \leq 1$.
 - For continuing tasks, $T = \infty$ and $0 < \gamma < 1$.

Bellman equation for a deterministic policy

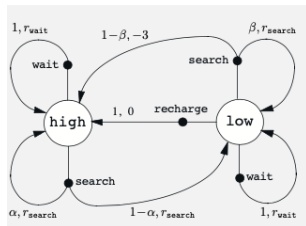
Backup diagram:



$$V^{\pi}(s) = \sum_{s' \in \mathcal{S}} p(s', s | \pi(a)) (r(s, s', \pi(a)) + \gamma V^{\pi}(s'))$$

For episodic tasks, set $V^{\pi}(s) = 0$ if s is a terminal state

Value function computation example



- Policy: $\pi(\text{high}) = \text{search}$, $\pi(\text{low}) = \text{search}$
- Bellman equations:

$$V^\pi(\text{h}) = \alpha(r_s + \gamma V^\pi(\text{h})) + (1 - \alpha)(r_s + \gamma V^\pi(1))$$

$$V^\pi(1) = (1 - \beta)(-3 + \gamma V^\pi(\text{h})) + \beta(r_s + \gamma V^\pi(1))$$

Solving the linear system

$$\begin{aligned} V^\pi(\mathbf{h}) &= \alpha(r_s + \gamma V^\pi(\mathbf{h})) + (1 - \alpha)(r_s + \gamma V^\pi(1)) \\ V^\pi(1) &= (1 - \beta)(-3 + \gamma V^\pi(\mathbf{h})) + \beta(r_s + \gamma V^\pi(1)) \end{aligned}$$

Rewrite in “linear algebra” style:

$$(I - \gamma P)V = b$$

Where:

$$P = \begin{bmatrix} \alpha & 1 - \alpha \\ 1 - \beta & \beta \end{bmatrix}, \quad \begin{bmatrix} r_s \\ -3(1 - \beta) + r_s\beta \end{bmatrix}$$

Numerical solution

Parameters:

$$\alpha = 0.8, \quad \beta = 0.3$$

$$r_{\text{search}} = 15, \quad r_{\text{wait}} = 10, \quad r_{\text{empty}} = -3.0$$

For $\gamma = 0.9$:

$$A = I - \gamma P = \begin{bmatrix} 0.28 & -0.18 \\ -0.63 & 0.73 \end{bmatrix}, \quad b = \begin{bmatrix} 15. \\ 2.4 \end{bmatrix}$$

Solution:

$$\begin{bmatrix} V(\text{high}) \\ V(\text{low}) \end{bmatrix} = \begin{bmatrix} 125.07692308 \\ 111.23076923 \end{bmatrix}$$

Solution strategies for a linear system

$$V = b + \gamma PV$$

- **Linear algebra solution.** Write the system in the form:

$$(I - \gamma P)V = b \quad \rightarrow \quad V = (I - \gamma P)^{-1}b$$

Don't use inversion, use an LU decomposition (Gaussian elimination).

- **Gauss-Jacobi Iteration.**

- Choose V_0 arbitrarily.
- Iterate:

$$V_{n+1} = b + \gamma PV_n$$

- This is an example of **fixed point iteration**.

Jacobi iteration for the RR model

Initial approximation: [0. 0.]

Step 1: [15. 2.4]

Step 2: [26.232 12.498]

Step 3: [36.13668 22.30062]

Step 4: [45.0325212 31.1872758]

Step 5: [53.03712491 39.19105282]

.....

Step 99: [125.07332253 111.22716869]

Step 100: [125.07368259 111.22752874]

Gauss-Seidel Iteration

- Similar to Jacobi iteration, but *use the already updated values in each step*:
 - Initialize V_0 arbitrarily.
 - Iterate:

$$\begin{aligned}V_{n+1}(1) &= b_1 + \gamma[p_{11}V_n(1) + p_{12}V_n(2)] \\V_{n+1}(2) &= b_2 + \gamma[p_{21}V_{n+1}(1) + p_{22}V_n(2)]\end{aligned}$$

- Results in slightly faster convergence.

Gauss-Seidel iteration for the RR model

Initial approximation: [0. 0.]

Step 1: [15. 11.85]

Step 2: [27.933 23.19729]

Step 3: [39.2872722 33.41424979]

Step 4: [49.30140095 42.48173004]

Step 5: [58.14372009 50.50061077]

.....

Step 79: [125.07002966 111.22451454]

Step 80: [125.07083397 111.22524433]

Optimal Value Function

- The optimal value function is:

$$V^*(s) = \max_{\pi} V^{\pi}(s) = \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t R_{t+1} \mid S_0 = s \right]$$

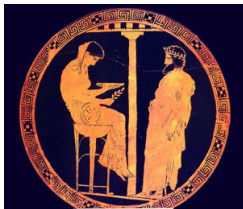
where π ranges over all admissible (randomized) policies.

- A policy π^* is optimal if:

$$V^{\pi^*}(s) = V^*(s) \quad \text{for all } s \in \mathcal{S}$$

- Under the finiteness assumptions we made, an optimal policy always exists.

How to live the best life

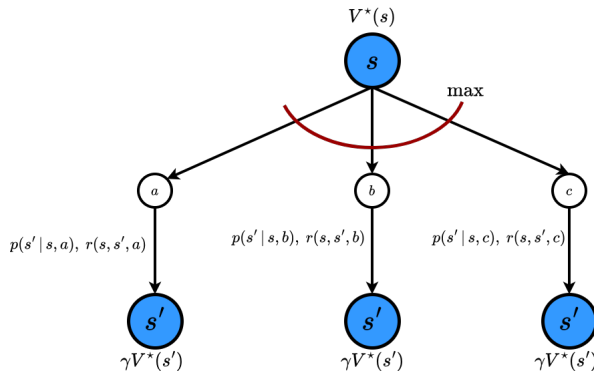


- There is an oracle that can tell us what to do to maximize our happiness. However, the oracle will only tell us what is the best we can do starting tomorrow.
- We can ask a series of questions to the oracle:
 - If I take action a_j today, what is the best outcome I can expect from tomorrow onward?
- We can then find:

$$\max\{R(a_j) + \mathbb{E}[\text{rewards from tomorrow onward} \mid \text{action } a_j \text{ chosen}]\}$$

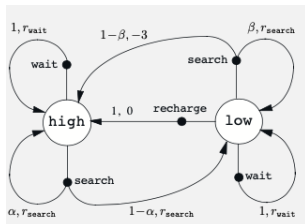
Optimal Bellman equation

Backup diagram:



$$V^*(s) = \max_{a \in \mathcal{A}(s)} \left\{ \sum_{s' \in \mathcal{S}} p(s, s' | a) [r(s, s', a) + \gamma V^*(s')] \right\}$$

Example — Recycling Robot



$$V^*(\text{high}) = \max \left\{ \alpha(r_{\text{search}} + \gamma V^*(\text{high})) + (1 - \alpha)(r_{\text{search}} + \gamma V^*(\text{low})), \right. \\ \left. r_{\text{wait}} + \gamma V^*(\text{high}) \right\}$$

$$V^*(\text{low}) = \max \left\{ (1 - \beta)(-3 + \gamma V^*(\text{high})) + \beta(r_{\text{search}} + \gamma V^*(\text{low})), \right. \\ \left. r_{\text{wait}} + \gamma V^*(\text{low}), \right. \\ \left. 0 + \gamma V^*(\text{high}) \right\}$$

Verification Theorem

Theorem

Suppose V^* is a solution of the optimal Bellman equation:

$$V^*(s) = \max_{a \in \mathcal{A}(s)} \left\{ \sum_{s' \in \mathcal{S}} p(s, s' | a) [r(s, s', a) + \gamma V^*(s')] \right\}.$$

Then, for any randomized policy $\pi(a | s)$ we have:

$$V^*(s) \geq V^\pi(s) \quad \text{for all } s \in \mathcal{S}$$

Furthermore, an optimal (deterministic) policy is given by:

$$\pi^*(s) = \arg \max \left\{ \sum_{s' \in \mathcal{S}} p(s, s' | a) [r(s, s', a) + \gamma V^*(s')] \mid a \in \mathcal{A}(s) \right\}$$

What are tabular solution methods

- In **tabular solution methods**, we build *complete* tables of the optimal value function and optimal policy.
- These methods can only be used to solve relatively small problems.
- It is also assumed that we have complete information about transition probabilities and rewards (*model based RL*).
- It is important to learn them, because:
 - The theory of these methods is *very* well established, so we know very well how they behave.
 - They are essential to develop intuition.
 - They contain the main ideas used in more advanced algorithms.
 - In practice, they provide a good way to build test cases for complex algorithms.

Value Iteration algorithm

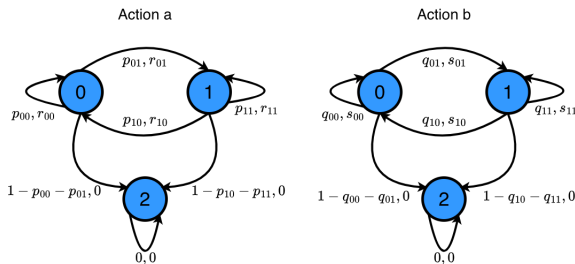
- Let $V_0(s)$ be arbitrary for $s \in \mathcal{S}$.
- Iterate:

$$V_{n+1}(s) = \max_{a \in \mathcal{A}(s)} \left\{ \sum_{s' \in \mathcal{S}} p(s, s' | a) [r(s, s', a) + \gamma V_n(s')] \right\},$$

except that we let $V_{n+1}(s) = 0$ if s is a terminal state in episodic tasks.

- Assume that:
 - $0 < \gamma < 1$ for continuing tasks.
 - $0 < \gamma \leq 1$ and a terminal state is reached with probability 1 for any policy.
- Then $\lim_{n \rightarrow \infty} V_n(s) = V^*(s)$ for all $s \in \mathcal{S}$.

Value Iteration Example — Tri-state MDP



Optimal Bellman equation for $\gamma = 1$:

$$V^*(0) = \max\{p_{00}(r_{00} + V^*(0)) + p_{01}(r_{01} + V^*(1)), \\ q_{00}(s_{00} + V^*(0)) + q_{01}(s_{01} + V^*(1))\}$$

$$V^*(1) = \max\{p_{10}(r_{10} + V^*(0)) + p_{11}(r_{11} + V^*(1)), \\ q_{10}(s_{10} + V^*(0)) + q_{11}(s_{11} + V^*(1))\}$$

Tri-state MDP — Value Iteration

$$V_{n+1}(0) = \max\{p_{00}(r_{00} + V_n(0)) + p_{01}(r_{01} + V_n(1)), \\ q_{00}(s_{00} + V_n(0)) + q_{01}(s_{01} + V_n(1))\}$$

$$V_{n+1}(1) = \max\{p_{10}(r_{10} + V_n(0)) + p_{11}(r_{11} + V_n(1)), \\ q_{10}(s_{10} + V_n(0)) + q_{11}(s_{11} + V_n(1))\}$$

Parameters for testing:

$$p_{00} = 0.2 \quad p_{01} = 0.7 \quad p_{10} = 0.5 \quad p_{11} = 0.3$$

$$q_{00} = 0.1 \quad q_{01} = 0.6 \quad q_{10} = 0.4 \quad q_{11} = 0.3$$

$$r_{00} = 10 \quad r_{01} = 15 \quad r_{10} = 8 \quad r_{11} = 12$$

$$s_{00} = 13 \quad s_{01} = 13 \quad s_{10} = 15 \quad s_{11} = 20$$

Value Iteration — Results

n	$V_n(0)$	$V_n(1)$
0	0.00	0.00
1	12.50	16.00
2	26.20	25.80
\vdots	\vdots	\vdots
38	71.24	63.57
39	71.25	63.57
40	71.25	63.57

- Convergence is not very fast.
- Optimal policy is found by a one-step lookahead. In this case we get:

$$\pi(0) = a, \quad \pi(1) = b$$

Gauss-Seidel Iteration

$$V_{n+1}(0) = \max\{p_{00}(r_{00} + V_n(0)) + p_{01}(r_{01} + V_n(1)), \\ q_{00}(s_{00} + V_n(0)) + q_{01}(s_{01} + V_n(1))\}$$

$$V_{n+1}(1) = \max\{p_{10}(r_{10} + V_{n+1}(0)) + p_{11}(r_{11} + V_n(1)), \\ q_{10}(s_{10} + V_{n+1}(0)) + q_{11}(s_{11} + V_n(1))\}$$

n	$V_n(0)$	$V_n(1)$
0	0.00	0.00
1	12.50	21.00
2	29.70	34.18
\vdots	\vdots	\vdots
26	71.24	63.57
27	71.25	63.57
28	71.25	63.57

- Easier to code, uses less memory.
- Slightly faster convergence.