

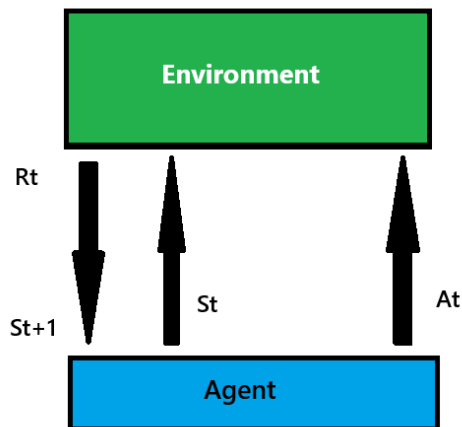
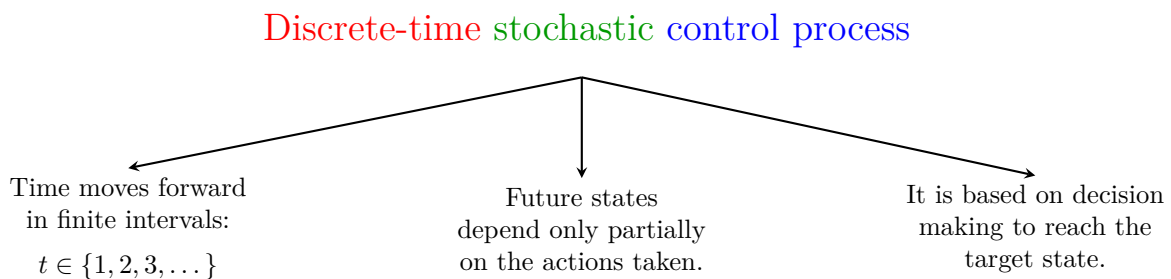
Reinforcement Learning

Reinforcement Learning is the branch of AI focused on solving control tasks, i.e. tasks that we have to perform with right actions in the right situation to get the best possible outcome.

1 Elements of Control Tasks

- State (S_t) : It is the relevant information that describes the situation the task environment is currently in.
- Action (A_t) : They are moves that the player can perform at any moment of time
- Reward (R_t) : It is a numerical value that the agent receives after carrying out an action.
- Agent : It is the entity that will participate in the task by observing its state and carrying out actions at each moment in time
- Environment : It comprises all of the aspects of the task that the agent doesn't control 100%.

2 Markov Decision Process



In many sequential decision-making tasks, time progresses in discrete steps — that is, $t = 1, 2, 3, \dots$. These problems can be effectively modeled using a framework known as a **Markov Decision Process (MDP)**.

An MDP describes the interaction between an **agent** and an **environment** over a sequence of time steps. The process unfolds as follows:

- At time $t = 0$, the agent observes the initial state of the environment.
- Based on this state, it selects an **action**.
- The action is executed, which causes the environment to transition to a new state.
- In response to this transition, the agent:
 - observes the new state,
 - and receives a scalar **reward signal**, indicating the quality or value of the transition.
- This cycle of *observe* \rightarrow *act* \rightarrow *transition* \rightarrow *reward* continues at each time step.

This discrete-time loop forms the basis of reinforcement learning problems, where the goal is to learn a policy that maximizes the cumulative reward over time.

An important advantage of the Markov Decision Process (MDP) framework is that it allows us to describe a control task in a compact and well-defined manner using just four core components:

- **State Space (\mathcal{S}):**
The set of all possible states in which the environment (or task) can exist. Each state represents a snapshot of the environment at a given time.
- **Action Space (\mathcal{A}):**
The set of all possible actions the agent can take while interacting with the environment.
- **Reward Function (\mathcal{R}):**
A mapping that assigns a numerical reward to each state-action pair (or state-action-next-state triplet). It quantifies the immediate benefit of taking a specific action in a given state.
- **Transition Probability Function (\mathcal{P}):**
Defines the probability of transitioning from one state to another, given a specific action. That is,

$$\mathcal{P}(s' \mid s, a) = \Pr(s_{t+1} = s' \mid s_t = s, a_t = a)$$

Together, these four elements — $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P})$ — fully specify a Markov Decision Process and allow us to formally describe and solve any control task.

The process has no memory

$$\mathcal{P}(s_{t+1} \mid s_t = s_t)$$

The next state depends only on the current state and not on the previous ones.

Why this is good:

- Simplifies modeling and computation
- Enables efficient algorithms (e.g., Q-learning, Dynamic Programming)
- Requires less memory — only the current state is needed

Why this can be bad:

- Ignores useful historical information
- Fails in partially observable environments

In practice, memory can be added using recurrent networks (e.g., LSTM) or by redefining the state to include sufficient history.

2.0.1 Transition Probability Function

Notation:

$$\mathcal{P}(s' | s, a) = \Pr(s_{t+1} = s' | s_t = s, a_t = a)$$

Meaning:

This represents the probability that the environment will transition to state s' at the next time step $t + 1$, given that it is currently in state s and the agent takes action a .

Intuition:

- Suppose you are controlling a robot.
 - State s = “robot is at location A”
 - Action a = “move forward”
- Possible next states s' :
 - With probability 0.9, the robot moves to location B (successful move)
 - With probability 0.1, the robot slips and stays at A

Then we have:

$$\mathcal{P}(B | A, \text{move}) = 0.9, \quad \mathcal{P}(A | A, \text{move}) = 0.1$$

These probabilities capture the uncertainty in the environment’s response to the agent’s actions.

2.1 Types of Markov Decision Processes

1. Finite vs. Infinite MDPs

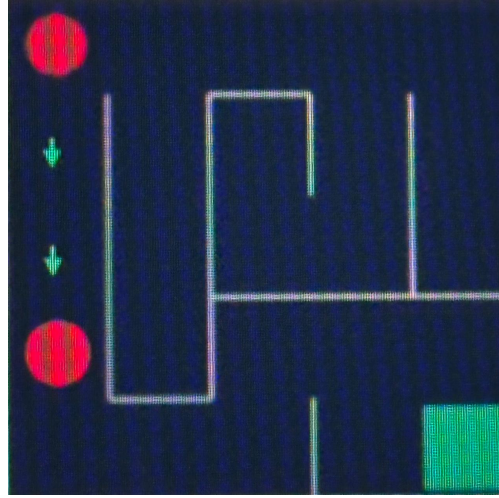
- **Finite MDP:** The sets of states, actions, and rewards are all finite.
Example: A 5×5 maze is a finite MDP:
 - 25 possible positions (states)
 - 4 possible actions per state (e.g., up, down, left, right)
 - Finite set of rewards
- **Infinite MDP:** At least one of the sets (states, actions, or rewards) is infinite, often due to continuous variables.
Example: Driving a car:
 - State includes position and speed (both continuous)
 - Infinite possible values (e.g., speed can be 75.0, 75.1, 75.01, etc.)

2. Episodic vs. Continuing MDPs

- **Episodic MDP:** The task terminates after a finite sequence of steps.
Example: Chess — the episode ends when the game concludes (e.g., checkmate).
- **Continuing MDP:** The task continues indefinitely without a terminal state.
Example: Monitoring or controlling an ongoing process (e.g., factory automation).

2.2 Trajectory and Episodes

- **Trajectory:** Elements that are generated when the agent moves from one state to another.



$$\tau = S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3$$

- **Episode:** Trajectory from the initial state of the task to a terminal state.

$$\tau = S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3 \rightarrow \text{Terminal}$$

In summary: All episodes are trajectories, but not all trajectories are episodes. An episode must terminate; a trajectory may be partial or infinite.

2.3 Reward vs Return

The goal of the task is represented by the rewards (R_t)

→ We want to maximize the sum of rewards

→ A short term reward can worsen long term results

Reward = R_t

Return, $G_t = R_{t+1} + R_{t+2} + \dots + R_t$

We want to maximize the episode's return

2.4 Discount Factor

The agent has no incentive to go to the goal through the shortest route

$$G_t = R_1 + R_2 + \dots + R_t$$

→ We have to modify somehow the return in a way that it rewards the most efficient behaviour.

We will modify future rewards by a discount factor γ

$$G_t = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots + \gamma^{T-t-1} R_t$$

$$\gamma \in [0, 1]$$

We want to maximize the long term sum of discounted rewards

Example:

Suppose an agent is navigating a grid world and receives the following rewards at each time step:

$$R_1 = 5, \quad R_2 = 3, \quad R_3 = 10$$

Let's assume the discount factor $\gamma = 0.9$. The discounted return G_t at time step $t = 1$ would be calculated as:

$$G_1 = R_1 = 5$$

At time step $t = 2$:

$$G_2 = R_1 + \gamma R_2 = 5 + 0.9 \times 3 = 5 + 2.7 = 7.7$$

At time step $t = 3$:

$$G_3 = R_1 + \gamma R_2 + \gamma^2 R_3 = 5 + 0.9 \times 3 + 0.9^2 \times 10 = 5 + 2.7 + 8.1 = 15.8$$

Thus, the agent's total discounted reward at each time step would be:

$$G_1 = 5, \quad G_2 = 7.7, \quad G_3 = 15.8$$

This example shows how future rewards are less valued, making the agent prefer immediate rewards or rewards sooner in time.

2.5 Policy

Policy is a function that decides what action to take in a particular state.

$$\pi : \mathcal{S} \rightarrow \mathcal{A}$$

Probability of taking action a in state s :

$$\pi(\mathcal{A}|\mathcal{S})$$

Action a taken in state s :

$$\pi(\mathcal{S})$$

Deterministic Policy:

$$\pi(s) = a$$

In a deterministic policy, the agent always chooses the same action for a given state.

Example: If the agent is in state $s = \text{"intersection"}$, then:

$$\pi(s) = \text{"turn left"}$$

This means the agent will always turn left when it reaches that intersection.

Stochastic Policy:

$$\pi(s) = [p(a_1), p(a_2), \dots, p(a_n)]$$

In a stochastic policy, the agent chooses an action based on a probability distribution over actions.

Example: In state $s = \text{"intersection"}$, the policy might be:

$$\pi(s) = [\text{turn left} : 0.6, \text{turn right} : 0.3, \text{go straight} : 0.1]$$

Here, the agent mostly prefers to turn left, but occasionally chooses other options based on the probabilities.

Summary:

- Deterministic policies are simpler and more predictable.
- Stochastic policies are more flexible and useful in uncertain or dynamic environments.

2.6 State Value Function

The **state value function**, denoted as $V^\pi(s)$, represents the expected return (cumulative reward) that an agent can obtain starting from state s and following policy π thereafter.

Definition:

$$V^\pi(s) = E_\pi [G_t \mid S_t = s]$$

where:

- $V^\pi(s)$: value of state s under policy π
- G_t : return from time t onward
- E_π : expectation over the trajectory induced by policy π

Expanded Equation:

$$V^\pi(s) = E_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

where $\gamma \in [0, 1]$ is the discount factor.

Example:

Consider a simple grid world where:

- The agent receives a reward of +1 for reaching the goal.
- From state $s = A$, it takes 3 steps on average to reach the goal.
- The policy π always chooses the shortest path.
- The discount factor $\gamma = 0.9$

Then the expected return from state A is:

$$V^\pi(A) = \gamma^0 \cdot 0 + \gamma^1 \cdot 0 + \gamma^2 \cdot 0 + \gamma^3 \cdot 1 = 0.9^3 = 0.729$$

So, under this policy, the value of state A is 0.729.

The value function helps the agent understand which states are more desirable in the long run.

Difference Between Return and State Value

Consider a simple environment in which an agent starts at state A and navigates step by step through different states to reach a goal state, where it receives a reward of +10.

Return: What Actually Happened in One Episode

Suppose the agent starts in state A at time $t = 0$ and reaches the goal in 3 steps. The rewards received during this episode are:

$$R_1 = 0, \quad R_2 = 0, \quad R_3 = 10$$

The return from time step 0 is then:

$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 = 0 + 0 + \gamma^2 \cdot 10 = \gamma^2 \cdot 10$$

This value G_0 represents the total reward collected in this single episode.

State Value: Expected Return from a State

Now suppose the agent plays 100 episodes, each starting from state A , following the same policy π . Due to randomness in the environment, the number of steps to reach the goal varies across episodes — some reach the goal in 3 steps, others in 5 steps, and some may not reach it at all. The return G_0 will vary in each episode. The **state value function** $V^\pi(A)$ is defined as the expected return starting from state A and following policy π :

$$V^\pi(A) = E_\pi[G_0 \mid S_0 = A]$$

In this case, $V^\pi(A)$ would be the average of the 100 returns observed, and it quantifies how valuable state A is under policy π in terms of long-term expected reward.

Conclusion:

The *return* is the actual reward accumulated in one specific run, while the *state value* is the expected return across many possible runs starting from the same state under a given policy.

2.7 Action-Value Function (Q-value)

The **action-value function**, also known as the **Q-value**, is denoted as $Q^\pi(s, a)$. It represents the expected return (cumulative discounted reward) starting from state s , taking action a , and then following policy π thereafter.

Definition:

$$Q^\pi(s, a) = E_\pi[G_t \mid S_t = s, A_t = a]$$

where:

- $Q^\pi(s, a)$ is the value of taking action a in state s under policy π
- G_t is the return from time step t onward
- The expectation is over all possible future state-action-reward sequences

Expanded Equation:

$$Q^\pi(s, a) = E_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s, A_t = a]$$

$$Q^\pi(s, a) = E_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a\right]$$

Interpretation:

While the state-value function $V^\pi(s)$ tells us how good it is to be in a state s , the Q-value $Q^\pi(s, a)$ tells us how good it is to *take action* a in state s , and then follow policy π .

Relation to Value Function:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a \mid s) \cdot Q^\pi(s, a)$$

This shows that the state-value is the expected Q-value, averaged over the policy's action probabilities.

Example:

Assume an agent in state $s = A$ has two actions:

- $Q^\pi(A, \text{"left"}) = 3.2$
- $Q^\pi(A, \text{"right"}) = 5.0$

Then, under policy π , the agent is more likely to choose the "right" action, as it has a higher expected return.

subsection Bellman Equations

The Bellman equations describe the recursive relationships for the value functions under a given policy π . They break down the expected return into immediate reward plus the discounted value of the next state.

1. Bellman Equation for State-Value Function $V^\pi(s)$:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} \mathcal{P}(s' | s, a) [\mathcal{R}(s, a, s') + \gamma V^\pi(s')]$$

- $\pi(a | s)$: probability of taking action a in state s under policy π
- $\mathcal{P}(s' | s, a)$: probability of transitioning to state s' given (s, a)
- $\mathcal{R}(s, a, s')$: expected immediate reward for that transition
- γ : discount factor

2. Bellman Equation for Action-Value Function $Q^\pi(s, a)$:

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} \mathcal{P}(s' | s, a) \left[\mathcal{R}(s, a, s') + \gamma \sum_{a' \in \mathcal{A}} \pi(a' | s') Q^\pi(s', a') \right]$$

- $Q^\pi(s, a)$: expected return of taking action a in state s and following π afterward
- The inner sum over a' reflects the expected value over the policy's actions at the next state

These recursive definitions are fundamental to many reinforcement learning algorithms, including policy evaluation, Q-learning, and value iteration.

2.8 Solving a Markov Decision Process (MDP)

The goal of solving an MDP is to find an **optimal policy** π^* that maximizes the expected return from each state.

2.8.1 Steps to Solve an MDP

1. Define the MDP:

- State space \mathcal{S}
- Action space \mathcal{A}
- Transition probability $\mathcal{P}(s' | s, a)$
- Reward function $\mathcal{R}(s, a, s')$
- Discount factor $\gamma \in [0, 1]$

2. Initialize the value function:

$$V_0(s) = 0 \quad \forall s \in \mathcal{S}$$

3. Iteratively apply the Bellman optimality equation:

$$V_{k+1}(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{P}(s' | s, a) [\mathcal{R}(s, a, s') + \gamma V_k(s')]$$

4. **Extract the optimal policy:**

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} \sum_{s'} \mathcal{P}(s' | s, a) [\mathcal{R}(s, a, s') + \gamma V^*(s')]$$

Example:

State B has two possible paths to reach terminal state D :

- **Path 1 (Longer path):** $B \xrightarrow{+1} s_1 \xrightarrow{+2} s_2 \xrightarrow{+10} D$
- **Path 2 (Shortcut):** $B \xrightarrow{+5} D$

Assume $\gamma = 0.9$. Calculate the value $V(B)$ using value iteration, and determine which path (i.e., which action) the optimal policy $\pi^*(B)$ will choose.

Step 0: Initialization

$$V_0(B) = V_0(s_1) = V_0(s_2) = V_0(D) = 0$$

Step 1: Evaluate terminal transitions

$$\begin{aligned} V_1(s_2) &= 10 + 0.9 \cdot 0 = 10 \\ V_1(s_1) &= 2 + 0.9 \cdot 10 = 11 \\ V_1(B_{\text{long}}) &= 1 + 0.9 \cdot 11 = 10.9 \\ V_1(B_{\text{short}}) &= 5 + 0.9 \cdot 0 = 5 \\ V_1(B) &= \max\{10.9, 5\} = 10.9 \end{aligned}$$

Step 2: Check for convergence

$$\begin{aligned} V_2(s_2) &= 10 \quad (\text{no change}) \\ V_2(s_1) &= 2 + 0.9 \cdot 10 = 11 \quad (\text{no change}) \\ V_2(B_{\text{long}}) &= 1 + 0.9 \cdot 11 = 10.9 \\ V_2(B_{\text{short}}) &= 5 \\ V_2(B) &= \max\{10.9, 5\} = 10.9 \end{aligned}$$

Final Answer:

$$V(B) = 10.9$$

Optimal Policy:

To find the best action from B , compare all available paths using:

$$\pi^*(B) = \arg \max_a [R(s, a, s') + \gamma V(s')]$$

This equation means: choose the action a that gives the highest total value.
In other words, take the action that gives the best immediate reward $R(s, a, s')$, plus the future value of the next state $V(s')$ (discounted by γ).

Since:

- Long path gives 10.9
- Shortcut gives 5

Therefore:

$$\pi^*(B) = \text{take the long path via } s_1 \rightarrow s_2 \rightarrow D$$

This path has a higher total discounted return, hence it is optimal.

3 Value Iteration

Value Iteration is a dynamic programming algorithm used to compute the optimal value function $V^*(s)$ in a Markov Decision Process (MDP). It iteratively updates the value of each state using the Bellman optimality equation:

$$V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')]$$

The process continues until the values converge, i.e., the change in value becomes smaller than a threshold ϵ .

Example (2x2 Grid):

$$\begin{bmatrix} A & B \\ C & \text{Terminal (+1)} \end{bmatrix}$$

Assume:

- $\gamma = 1$
- Terminal state at $[1][1]$ with reward +1 and value fixed at $V = 1$
- All other rewards = 0
- Transitions are deterministic (e.g., up/down/left/right always succeed)

Initialization:

$$V_0(A) = V_0(B) = V_0(C) = 0, \quad V_0(\text{Terminal}) = 1$$

Iteration 1:

- $V_1(C) = \max \{0 + \gamma \cdot V_0(\text{Terminal})\} = 1$
- $V_1(B)$ and $V_1(A)$ can only reach states with value 0 \rightarrow remain 0

Iteration 2:

- $V_2(B) = \max \{0 + \gamma \cdot V_1(\text{Terminal}), 0 + \gamma \cdot V_1(A)\} = 1$
- $V_2(A) = \max \{0 + \gamma \cdot V_1(B), 0 + \gamma \cdot V_1(C)\} = 1$

Final Value Table (after convergence):

$$V^*(A) = 1, \quad V^*(B) = 1, \quad V^*(C) = 1, \quad V^*(\text{Terminal}) = 1$$

All states learn to reach the terminal state to get the highest reward.

MDP Grid Environment

.	.	.	1
.	W	.	-1
.	.	.	.

Listing 1: Value Iteration in a Grid World

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

# Environment parameters
gamma = 0.9
noise = 0.2
living_reward = 0

# Grid dimensions
rows, cols = 3, 4

# Define the rewards grid
# None = wall, other values = rewards
rewards = np.zeros((rows, cols))
rewards[0, 3] = +1 # Terminal state with +1
rewards[1, 3] = -1 # Terminal state with -1
rewards[1, 1] = None # Wall

# Terminal and wall states
terminal_states = [(0, 3), (1, 3)]
wall_states = [(1, 1)]

rewards, terminal_states, wall_states

# Display the matrix
display_matrix = np.full((rows, cols), '', dtype=object)
for r in range(rows):
    for c in range(cols):
        if (r, c) in terminal_states:
            display_matrix[r, c] = f"{rewards[r, c]:.0f}"
        elif (r, c) in wall_states:
            display_matrix[r, c] = "W"
        else:
            display_matrix[r, c] = "."

fig, ax = plt.subplots(figsize=(6, 4))
ax.set_xticks(np.arange(cols + 1) - 0.5, minor=True)
ax.set_yticks(np.arange(rows + 1) - 0.5, minor=True)
ax.grid(which="minor", color="black", linestyle='--', linewidth=1)
ax.tick_params(which="minor", bottom=False, left=False)

for r in range(rows):
    for c in range(cols):
        text = display_matrix[r, c]
        ax.text(c, r, text, va='center', ha='center', fontsize=16)

ax.set_xlim(-0.5, cols - 0.5)
ax.set_ylim(rows - 0.5, -0.5)
ax.set_xticks([])
ax.set_yticks([])

plt.title("MDP Grid Environment")
plt.show()

# Actions and their directions
actions = ['U', 'D', 'L', 'R']
action_vectors = {
```

```

    'U': (-1, 0),
    'D': (1, 0),
    'L': (0, -1),
    'R': (0, 1)
}

# Initialize value function
V = np.zeros((rows, cols))
policy = np.full((rows, cols), '', dtype=object)

def in_bounds(state):
    r, c = state
    return 0 <= r < rows and 0 <= c < cols and (r, c) not in wall_states

def compute_action_value(state, action, V):
    r, c = state
    primary_move = action_vectors[action]

    if action in ['U', 'D']:
        sideways = ['L', 'R']
    else:
        sideways = ['U', 'D']

    moves = [(primary_move, 1 - noise)] + [(action_vectors[a], noise / 2) for a in
    ↪ sideways]

    value = 0
    for move, prob in moves:
        new_r, new_c = r + move[0], c + move[1]
        if not in_bounds((new_r, new_c)):
            new_r, new_c = r, c
        reward = living_reward if (new_r, new_c) not in terminal_states else
        ↪ rewards[new_r, new_c]
        value += prob * (reward + gamma * V[new_r, new_c])
    return value

theta = 1e-4
iteration = 0
while True:
    delta = 0
    new_V = np.copy(V)
    for r in range(rows):
        for c in range(cols):
            if (r, c) in terminal_states or (r, c) in wall_states:
                continue
            values = [compute_action_value((r, c), a, V) for a in actions]
            best_value = max(values)
            new_V[r, c] = best_value
            delta = max(delta, abs(best_value - V[r, c]))
    V = new_V
    iteration += 1
    if delta < theta:
        break

for r in range(rows):
    for c in range(cols):
        if (r, c) in terminal_states:
            policy[r, c] = 'T'
        elif (r, c) in wall_states:

```

```

        policy[r, c] = 'W'
    else:
        values = [compute_action_value((r, c), a, V) for a in actions]
        best_action = actions[np.argmax(values)]
        policy[r, c] = best_action

import pandas as pd
import ace_tools as tools; tools.display_dataframe_to_user(name="Value Function",
    ↪ dataframe=pd.DataFrame(V))

policy_display = pd.DataFrame(policy)
tools.display_dataframe_to_user(name="Policy", dataframe=policy_display)

```

	0	1	2	3
1	R	R	R	T
2	U	W	U	T
3	U	L	U	L

4 Policy Iteration

Policy Iteration is a dynamic programming algorithm used to compute the optimal policy π^* in a Markov Decision Process (MDP). It consists of two main steps that are repeated until convergence:

1. **Policy Evaluation:** For a fixed policy π , compute the value function $V^\pi(s)$ by solving:

$$V^\pi(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

2. **Policy Improvement:** Update the policy by choosing the action that yields the highest expected return:

$$\pi_{\text{new}}(s) = \arg \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^\pi(s')]$$

Repeat these two steps until the policy no longer changes. This results in the optimal value function V^* and the optimal policy π^* .

Key Advantage: Policy Iteration often converges faster than Value Iteration in practice, since it uses full value evaluation before each improvement step.

Listing 2: Policy Iteration in Grid World

```

import numpy as np

# Environment parameters
gamma = 0.9
noise = 0.2
living_reward = 0

# Grid dimensions
rows, cols = 3, 4

# Define rewards grid

```

```

rewards = np.zeros((rows, cols))
rewards[0, 3] = +1
rewards[1, 3] = -1
rewards[1, 1] = None # Wall

# Terminal and wall states
terminal_states = [(0, 3), (1, 3)]
wall_states = [(1, 1)]

# Actions and their vectors
actions = ['U', 'D', 'L', 'R']
action_vectors = {
    'U': (-1, 0),
    'D': (1, 0),
    'L': (0, -1),
    'R': (0, 1)
}

# Utility function to check bounds
def in_bounds(state):
    r, c = state
    return 0 <= r < rows and 0 <= c < cols and (r, c) not in wall_states

# Compute expected value of taking an action at a state
def compute_action_value(state, action, V):
    r, c = state
    primary_move = action_vectors[action]

    if action in ['U', 'D']:
        sideways = ['L', 'R']
    else:
        sideways = ['U', 'D']

    moves = [(primary_move, 1 - noise)] + [(action_vectors[a], noise / 2) for a in sideways]

    value = 0
    for move, prob in moves:
        new_r, new_c = r + move[0], c + move[1]
        if not in_bounds((new_r, new_c)):
            new_r, new_c = r, c
        reward = living_reward if (new_r, new_c) not in terminal_states else rewards[new_r, new_c]
        value += prob * (reward + gamma * V[new_r, new_c])
    return value

# Initialize policy randomly
policy = np.full((rows, cols), 'U', dtype=object)
for r, c in terminal_states + wall_states:
    policy[r, c] = None

# Initialize value function
V = np.zeros((rows, cols))

# Policy Iteration
is_policy_stable = False
iteration = 0

while not is_policy_stable:

```

```

# Policy Evaluation
while True:
    delta = 0
    new_V = np.copy(V)
    for r in range(rows):
        for c in range(cols):
            if (r, c) in terminal_states or (r, c) in wall_states:
                continue
            v = compute_action_value((r, c), policy[r, c], V)
            new_V[r, c] = v
            delta = max(delta, abs(v - V[r, c]))
    V = new_V
    if delta < 1e-4:
        break

# Policy Improvement
is_policy_stable = True
for r in range(rows):
    for c in range(cols):
        if (r, c) in terminal_states or (r, c) in wall_states:
            continue
        old_action = policy[r, c]
        action_values = [compute_action_value((r, c), a, V) for a in actions]
        best_action = actions[np.argmax(action_values)]
        policy[r, c] = best_action
        if best_action != old_action:
            is_policy_stable = False

import pandas as pd
import ace_tools as tools; tools.display_dataframe_to_user(name="Policy Iteration
    ↪ - Value Function", dataframe=pd.DataFrame(V))

policy_display = pd.DataFrame(policy)
tools.display_dataframe_to_user(name="Policy Iteration - Policy", dataframe=
    ↪ policy_display)

```

5 Monte Carlo Methods

Monte Carlo methods are a class of algorithms used in Reinforcement Learning (RL) to estimate value functions and learn optimal policies by averaging returns from sampled episodes. **Unlike dynamic programming methods, Monte Carlo methods do not require a complete model of the environment (i.e., transition probabilities or reward functions). Instead, they learn directly from experience.**

Why Monte Carlo?

- Can be used in model-free settings, where the environment's dynamics are unknown.
- Suitable for episodic tasks where learning happens at the end of each episode.
- Converges to accurate value estimates given enough samples (by the Law of Large Numbers).

How it Works:

1. Generate episodes by interacting with the environment using a policy π .
2. For each state (or state-action pair) visited, observe the return G_t after that visit.
3. Use the average of these returns to estimate the value function $V(s)$ or $Q(s, a)$.
4. Optionally, improve the policy using the estimated values (e.g., ϵ -greedy).

Two Variants:

- **First-visit MC:** Only the first occurrence of a state in each episode is used for updating the value.
- **Every-visit MC:** Every occurrence of the state in an episode is used for updating.

Value Estimation Equation:

$$V(s) \approx \frac{1}{N(s)} \sum_{i=1}^{N(s)} G^{(i)} \quad \text{where } G^{(i)} \text{ is the return after } i^{\text{th}} \text{ visit to } s$$

Monte Carlo methods are powerful when dealing with environments where learning is based purely on sampling, such as in games, simulations, or real-world robotics.