



EXPERIMENT ASSESSMENT

ACADEMIC YEAR 2023-24

Course: Reinforcement Learning Lab

Course code: CSDOL8013

Year: BE **Sem:** VIII

Experiment No.: 4
Aim: - Applying dynamic programming algorithms, such as policy evaluation and policy improvement, to solve a small-scale MDP problem.
Name:
Roll Number:
Date of Performance:
Date of Submission:

Evaluation

Performance Indicator	Max. Marks	Marks Obtained
Performance	5	
Understanding	5	
Journal work and timely submission.	10	
Total	20	

Performance Indicator	Exceed Expectations (EE)	Meet Expectations (ME)	Below Expectations (BE)
Performance	5	3	2
Understanding	5	3	2
Journal work and timely submission.	10	8	4

Checked by

Name of Faculty : Rujuta Vartak

Signature :

Date :



EXPERIMENT 4

Aim: Applying dynamic programming algorithms, such as policy evaluation and policy improvement, to solve a small-scale MDP problem.

Objective: The objective is to employ dynamic programming algorithms, such as policy evaluation and policy improvement, to effectively address a small-scale Markov Decision Process (MDP) problem, aiming to iteratively refine and optimize policies, thereby gaining insights into foundational concepts of reinforcement learning and optimization.

Theory:

Markov Decision Processes (MDPs): MDPs are mathematical frameworks used to model decision-making problems where outcomes are partly random and partly under the control of a decision-maker. They consist of states, actions, transition probabilities, rewards, and a discount factor.

Policy Evaluation: Policy evaluation is the process of determining the value function for a given policy. The value function represents the expected cumulative reward starting from a particular state and following a specific policy. The basic idea is to iteratively update the value of each state until convergence.

Algorithm for policy evaluation :

Input:

MDP: (S, A, P, R, γ) , where

S is the set of states.

A is the set of actions.

P is the transition probability matrix, $P(s' | s, a)$, representing the probability of transitioning to state s' from state s by taking action a .

R is the reward function, $R(s, a, s')$, representing the immediate reward received after transitioning from state s to state s' by taking action a .

γ is the discount factor.

Output:

Value function $V(s)$ for each state s .

Algorithm:

Initialize the value function arbitrarily: $V(s)$ for all s in S .

Repeat until convergence:



Initialize Δ to be a very large value (e.g., infinity).

For each state s in S :

Let v be the current value of $V(s)$.

Update $V(s)$ using the Bellman expectation equation:

$V(s) = \sum [P(s' | s, a) * (R(s, a, s') + \gamma * V(s'))]$ over all possible next states s' and actions a .

Update Δ to be the maximum between Δ and $|v - V(s)|$.

If Δ is smaller than a predefined threshold ϵ , break.

Return the converged value function $V(s)$.

Explanation:

The algorithm iteratively updates the value function $V(s)$ for each state s using the Bellman expectation equation until convergence.

At each iteration, it computes the expected value of being in state s and taking action a , which is the sum of immediate reward $R(s, a, s')$ and the discounted value of the next state $V(s')$.

The process continues until the maximum change in the value function between iterations (Δ) falls below a predefined threshold ϵ , indicating convergence.

The output is the converged value function $V(s)$, which represents the expected cumulative reward from each state under the given policy.

This algorithm is known as "Iterative Policy Evaluation" and is a fundamental component of dynamic programming approaches for solving MDPs. It provides a way to estimate the value of each state under a given policy.

Policy Improvement: Policy improvement involves selecting better actions in each state to improve the current policy. It's based on the idea of greedily selecting actions that maximize the expected cumulative reward given the current value function.

Algorithm for policy improvement:

Input:

MDP: (S, A, P, R, γ) , where

S is the set of states.

A is the set of actions.

P is the transition probability matrix, $P(s' | s, a)$, representing the probability of transitioning to state s' from state s by taking action a .

R is the reward function, $R(s, a, s')$, representing the immediate reward received after transitioning from state s to state s' by taking action a .

γ is the discount factor.

Value function $V(s)$ for each state s .

Output:



Vidyavardhini's College of Engineering & Technology

Department of Computer Science & Engineering (Data Science)

Improved policy $\pi'(s)$ for each state s .

Algorithm:

Initialize a boolean variable `policy_stable` to true.

For each state s in S , do:

Let `old_action` be the current action selected by the policy π for state s .

Compute Q-value for each action a in A :

$Q(s, a) = \sum [P(s' | s, a) * (R(s, a, s') + \gamma * V(s'))]$ over all possible next states s' .

Select the action `a_max` that maximizes the Q-value: `a_max = argmax(Q(s, a))`.

Update the policy $\pi'(s)$ to select the action `a_max`.

Check for policy stability:

If the new policy π' is different from the old policy π for any state s , set `policy_stable` to false.

If `policy_stable` is true, return the improved policy π' .

Else, return to step 2 and repeat the process with the updated policy π' .

Explanation:

The algorithm iterates over each state in the state space and computes the Q-value for each action based on the current value function $V(s)$.

It selects the action that maximizes the Q-value as the new action for the state.

The process continues until the policy stabilizes, i.e., the new policy is the same as the old policy for all states.

The output is the improved policy π' that greedily selects actions to maximize the expected cumulative reward according to the current value function $V(s)$.

This algorithm is known as "Policy Improvement" and is used in combination with policy evaluation to iteratively improve the policy until convergence to an optimal policy in dynamic programming approaches for solving MDPs.

1. **Iterative Policy Evaluation and Policy Improvement:** The process of policy evaluation and policy improvement is often interleaved. After evaluating a policy, we improve it by selecting better actions based on the updated value function. Then, we re-evaluate the policy to refine the value estimates further.
2. **Convergence:** Both policy evaluation and policy improvement converge to the optimal value function and policy if executed iteratively until convergence.
3. **Implementation:** Dynamic programming algorithms can be implemented efficiently using programming languages like Python, where you can define MDPs, transition probabilities, rewards, value functions, and policies, and then iteratively update them until convergence.

This general approach can be applied to small-scale MDP problems to find the optimal policy



efficiently. However, for larger problems, approximate methods like reinforcement learning techniques may be more practical due to the computational complexity of dynamic programming algorithms.

CODE

```
import numpy as np

# Define MDP parameters
NUM_STATES = 3
NUM_ACTIONS = 2
GAMMA = 0.9 # Discount factor

# Define transition probabilities
# transition_probs[state][action][next_state]
transition_probs = np.array([
    [[0.2, 0.4, 0.4], [1.0, 0.0, 0.0]], # from S1
    [[0.0, 0.0, 1.0], [0.0, 0.0, 1.0]], # from S2
    [[1.0, 0.0, 0.0], [1.0, 0.0, 0.0]] # from S3
])

# Define rewards
# rewards[state][action][next_state]
rewards = np.array([
    [[-1, -1, -1], [0, 0, 0]], # from S1
    [[-1, -1, -1], [-1, -1, 10]], # from S2
    [[-1, -1, -1], [-1, -1, -1]] # from S3
])

# Initialize a random policy
policy = np.random.randint(0, NUM_ACTIONS, size=NUM_STATES)

def policy_evaluation(policy, transition_probs, rewards, gamma=0.9,
tol=1e-6):
    V = np.zeros(NUM_STATES) # Initialize value function to zeros
    while True:
        delta = 0
        for s in range(NUM_STATES):
            v = V[s]
            bellman_expectation = sum(transition_probs[s][policy[s]][s1] *
(rewards[s][policy[s]][s1] + gamma * V[s1]) for s1 in range(NUM_STATES))
            V[s] = bellman_expectation
            delta = max(delta, abs(v - V[s]))
```



Vidyavardhini's College of Engineering & Technology

Department of Computer Science & Engineering (Data Science)

```
if delta < tol:
    break
return V

def policy_improvement(V, transition_probs, rewards, gamma=0.9):
    policy_stable = True
    for s in range(NUM_STATES):
        old_action = policy[s]
        action_values = np.zeros(NUM_ACTIONS)
        for a in range(NUM_ACTIONS):
            action_values[a] = sum(transition_probs[s][a][s1] *
(rewards[s][a][s1] + gamma * V[s1]) for s1 in range(NUM_STATES))
        # Greedily select the best action
        policy[s] = np.argmax(action_values)
        if old_action != policy[s]:
            policy_stable = False
    return policy, policy_stable

# Perform policy iteration
policy_stable = False
iteration = 0
while not policy_stable:
    print(f"Iteration {iteration}: Policy {policy}")
    V = policy_evaluation(policy, transition_probs, rewards, gamma=GAMMA)
    policy, policy_stable = policy_improvement(V, transition_probs, rewards,
gamma=GAMMA)
    iteration += 1

print(f"Optimal Policy: {policy}")
print(f"Optimal Value Function: {V}")
```

Output:

```
➡ Iteration 0: Policy [0 0 0]
Iteration 1: Policy [1 1 0]
Iteration 2: Policy [0 1 0]
Optimal Policy: [0 1 0]
Optimal Value Function: [ 9.37377142 16.69275405  7.43639428]
```



Conclusion:

1. Give one example of dynamic Programming.

One classic example of dynamic programming is the "Fibonacci sequence". The Fibonacci sequence is a series of numbers in which each number (except for the first two) is the sum of the two preceding ones. It starts with 0 and 1, and the sequence goes like this: 0, 1, 1, 2, 3, 5, 8, 13, 21, ... The naive recursive implementation of Fibonacci can be very inefficient because it recalculates the values for the same subproblems multiple times. Dynamic programming can optimize this by storing the results of subproblems and reusing them when needed. Instead of recalculating Fibonacci numbers from scratch for each value of n , we store previously calculated Fibonacci numbers in the fib list and reuse them as needed. This significantly improves the efficiency of the Fibonacci calculation, especially for large values of n . This approach is an example of bottom-up dynamic programming, where we solve smaller subproblems first and build up to the larger problem.

2. Explain how dynamic programming is utilized in reinforcement learning.

Here's how dynamic programming is utilized in reinforcement learning:

- a. Policy Evaluation: In RL, the value function represents the expected return (or cumulative reward) that an agent can achieve from a given state under a certain policy. Dynamic programming is used to iteratively estimate the value function for a given policy until convergence. This process is called policy evaluation.
- b. Policy Improvement: Once the value function has been evaluated, dynamic programming is used to improve the policy based on the estimated value function. The policy improvement step involves selecting actions that lead to states with higher value estimates.
- c. Value Iteration: Value iteration is another DP algorithm commonly used in reinforcement learning. It combines both policy evaluation and policy improvement into a single step. In value iteration, the value function is updated iteratively by taking the maximum expected return over all possible actions from each state.
- d. Model-Based RL: Dynamic programming methods are particularly useful in model-based RL, where the agent has access to a model of the environment dynamics (transition probabilities and rewards).