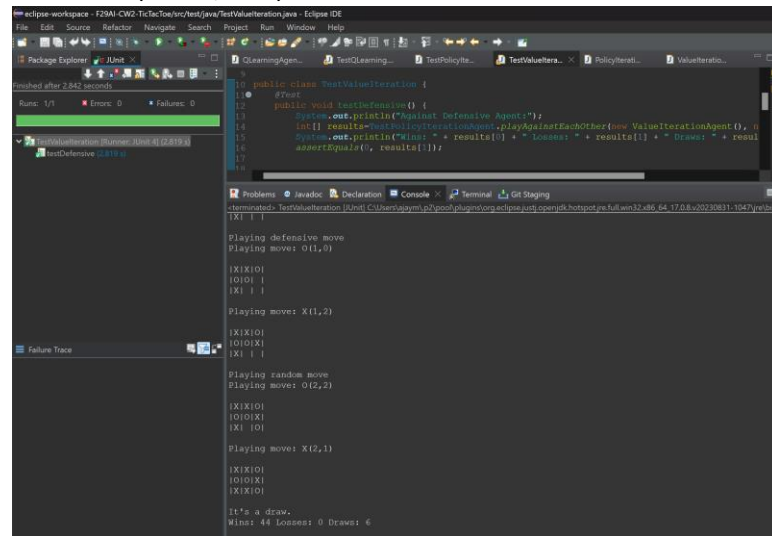


F29AI – CW2 2023 – 2024 VALUE ITERATION [AJAY MENON, H00418802]

Defensive (W: 44, D: 6)



```
10 public class TestValueIteration {
11     @Test
12     public void testDefensive() {
13         System.out.println("Against Defensive Agent");
14         int[] results = TestPolicyIterationAgent.playAgainstEachOther(new ValueIterationAgent(), 0);
15         System.out.println("Wins: " + results[0] + " Losses: " + results[1] + " Draws: " + results[2]);
16         assertEquals(0, results[1]);
17     }
18 }

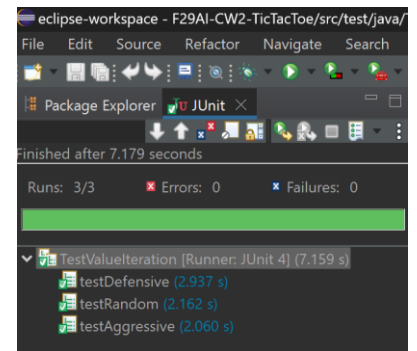
Playing defensive move
Playing move: O(1,0)
|X|X|O|
|O|O| |
|X| | |

Playing move: X(1,2)
|X|X|O|
|O|O|X|
|X| | |

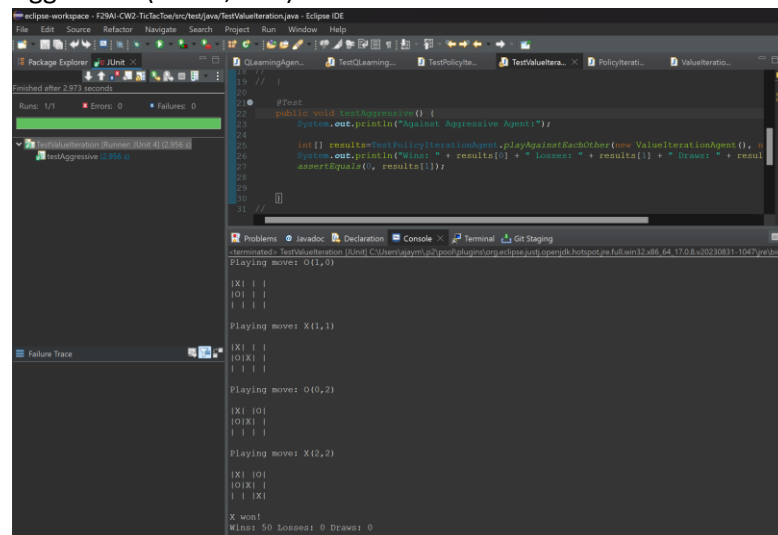
Playing random move
Playing move: O(2,2)
|X|X|O|
|O|O|X|
|X| |O|

Playing move: X(2,1)
|X|X|O|
|O|O|X|
|X|X|O|

It's a draw.
Wins: 44 Losses: 0 Draws: 6
```



Aggressive (W: 50, D: 0)



```
21 @Test
22 public void testAggressive() {
23     System.out.println("Against Aggressive Agent");
24     int[] results = TestPolicyIterationAgent.playAgainstEachOther(new ValueIterationAgent(), 0);
25     System.out.println("Wins: " + results[0] + " Losses: " + results[1] + " Draws: " + results[2]);
26     assertEquals(0, results[1]);
27 }
28 //
29
30 //

Playing move: O(1,0)
|X| | |
|O| | |
| | | |

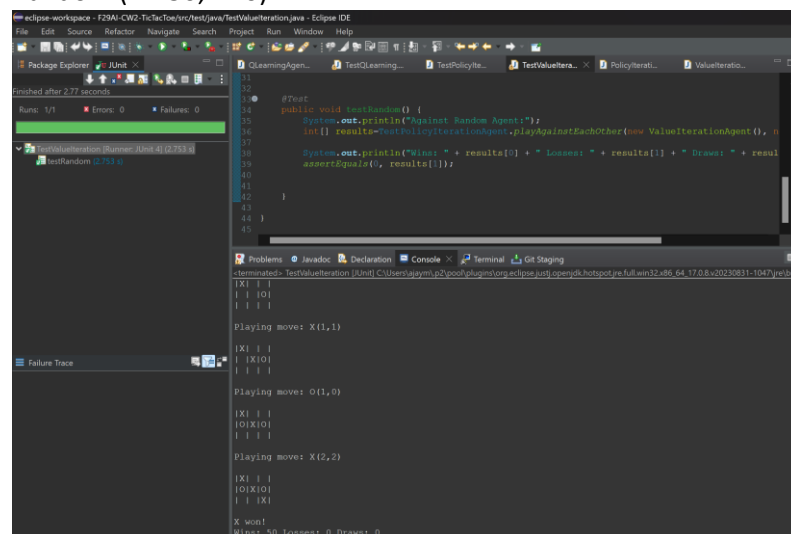
Playing move: X(1,1)
|X| | |
|O|X| |
| | | |

Playing move: O(0,2)
|X| |O|
|O|X| |
| | | |

Playing move: X(2,2)
|X| |O|
|O|X| |
| | |X|

X won!
Wins: 50 Losses: 0 Draws: 0
```

Random (W: 50, D: 0)



```
32 @Test
33 public void testRandom() {
34     System.out.println("Against Random Agent");
35     int[] results = TestPolicyIterationAgent.playAgainstEachOther(new ValueIterationAgent(), 0);
36     System.out.println("Wins: " + results[0] + " Losses: " + results[1] + " Draws: " + results[2]);
37     assertEquals(0, results[1]);
38 }

Playing move: X(1,1)
|X| | |
| | |O|
| | | |

Playing move: O(1,0)
|X| | |
|O|X|O|
| | | |

Playing move: X(2,2)
|X| | |
|O|X|O|
| | |X|

X won!
Wins: 50 Losses: 0 Draws: 0
```

ChatGPT Conversation SS:

(Chat Link: <https://chat.openai.com/share/ea1bda07-23cc-43d7-83d6-1400ba1d168f>)

You

Explain value iteration in terms of Artificial Intelligence



ChatGPT

Value iteration is a fundamental concept in the field of artificial intelligence, particularly in the domain of reinforcement learning and dynamic programming. It is used to find the optimal value function for a Markov Decision Process (MDP), which is a mathematical model for decision-making in situations where the outcome is uncertain.

Here's a breakdown of value iteration in the context of artificial intelligence:

1. **Markov Decision Process (MDP):** An MDP is a mathematical model that represents decision-making in a situation where the outcome is uncertain. It consists of a set of states, actions, transition probabilities, rewards, and discount factor.
2. **Value Function:** The value function, denoted as $V(s)$, represents the expected cumulative reward that an agent can obtain from a given state s by following a particular policy. The goal of value iteration is to compute the optimal value function, which provides the maximum expected cumulative reward for each state.
3. **Bellman Equation:** The Bellman equation expresses the relationship between the value of a state and the values of its neighboring states. For a given state s , the Bellman equation is defined as:
$$V(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s')]$$

Here,

 - $R(s, a)$ is the immediate reward of taking action a in state s .
 - $P(s'|s, a)$ is the probability of transitioning to state s' from state s after taking action a .
 - γ is the discount factor, which represents the agent's preference for present rewards over future rewards.
4. **Value Iteration Algorithm:**
 - Initialize the value function $V(s)$ arbitrarily for all states.
 - Repeat until convergence:
 - For each state s , update the value function using the Bellman equation.
$$V(s) \leftarrow \max_{a \in A} [R(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s')]$$
 - Convergence occurs when the change in the value function becomes negligible.
5. **Optimal Policy Extraction:** Once the value iteration process converges, the optimal policy can be extracted by choosing actions that maximize the expression within the square brackets in the Bellman equation for each state.
6. **Convergence:** Value iteration converges to the optimal value function, and consequently, the optimal policy. The convergence is guaranteed due to the monotonicity property of the Bellman operator.

In summary, value iteration is an iterative algorithm used in AI to find the optimal value function for decision-making under uncertainty. It is a key component of reinforcement learning and dynamic programming approaches to solving MDPs.



You

Provide a suitable pseudocode for value iteration



ChatGPT

Certainly! Here's a simple pseudocode for the Value Iteration algorithm:

plaintext

Copy code

```
function ValueIteration(MDP):
    // MDP is a tuple (S, A, P, R, γ), where:
    // S is the set of states
    // A is the set of actions
    // P is the transition probability function
    // R is the reward function
    // γ is the discount factor

    // Initialization
    for each state s in S:
        V[s] = 0

    // Value Iteration
    repeat:
        Δ = 0
        for each state s in S:
            v = V[s]
            V[s] = max_a Σ_{s'} P(s' | s, a) * [R(s, a, s') + γ * V[s']]
            Δ = max(Δ, |v - V[s]|)
        until Δ < ε // ε is a small positive threshold for convergence

    // Policy Extraction
    for each state s in S:
        π[s] = argmax_a Σ_{s'} P(s' | s, a) * [R(s, a, s') + γ * V[s']]

    return V, π
```

In this pseudocode:

- S is the set of states,
- A is the set of actions,
- P is the transition probability function,
- R is the reward function,
- γ is the discount factor,
- V is the array representing the value function for each state,
- π is the array representing the optimal policy for each state,
- Δ is the change in the value function, and
- ϵ is a small positive threshold for convergence.

The algorithm iteratively updates the value function until it converges, and then extracts the optimal policy based on the computed values. The convergence criterion is typically based on the change in the value function being below a certain threshold (ϵ).

📄 👍 🗨️ ↻