



# Introduction to Reinforcement Learning

## Dynamic Programming - Value Iteration

Sarwan Ali

Department of Computer Science  
Georgia State University

 Dynamic Programming Methods 

# Today's Learning Journey

- 1 Review: Bellman Equations
- 2 Value Iteration Algorithm
- 3 Convergence Theory
- 4 Implementation Details
- 5 Worked Example
- 6 Computational Complexity
- 7 Advantages and Limitations
- 8 Comparison with Policy Iteration
- 9 Extensions and Applications
- 10 Summary

# Bellman Equations Recap

## State Value Function:

$$v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] \quad (1)$$

$$= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad (2)$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \quad (3)$$

## Optimal State Value Function:

$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \quad (4)$$

$$= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \quad (5)$$

## Key Insight

The Bellman optimality equation provides the foundation for value iteration

# Value Iteration: Core Idea

## Principle

**Value Iteration** turns the Bellman optimality equation into an iterative update rule

**Instead of solving:**  $v_*(s) = \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma v_*(s')]$

**We iterate:**  $v_{k+1}(s) = \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma v_k(s')]$

## Key Properties:

- Model-based approach
- Guaranteed convergence
- Finds optimal policy

## Requirements:

- Complete MDP model
- Finite state/action spaces
- Discount factor  $\gamma < 1$

# Value Iteration Algorithm

---

## Algorithm 1 Value Iteration

---

**Require:** MDP  $(S, A, P, R, \gamma)$ , threshold  $\theta > 0$

- 1: Initialize  $V(s) = 0$  for all  $s \in S$
  - 2: **repeat**
  - 3:    $\Delta \leftarrow 0$
  - 4:   **for** each state  $s \in S$  **do**
  - 5:      $v \leftarrow V(s)$
  - 6:      $V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$
  - 7:      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
  - 8:   **end for**
  - 9: **until**  $\Delta < \theta$
  - 10: **Output:** Optimal policy  $\pi_*(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$
- 

## Time Complexity

$O(|S|^2|A|)$  per iteration, polynomial convergence

# Value Iteration: Step-by-Step Breakdown

## Step 1: Initialization

- Set  $V_0(s) = 0$  for all states  $s$
- Choose convergence threshold  $\theta$

## Step 2: Value Update (Bellman Backup)

$$V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')] \quad (6)$$

## Step 3: Convergence Check

- Compute  $\Delta = \max_s |V_{k+1}(s) - V_k(s)|$
- Stop if  $\Delta < \theta$

## Step 4: Policy Extraction

$$\pi_*(s) = \arg \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_*(s')] \quad (7)$$

# Convergence Guarantees

## Theorem (Value Iteration Convergence)

*For any finite MDP with  $\gamma < 1$ , value iteration converges to the unique optimal value function  $v_*$ .*

### Proof Sketch:

- 1 The Bellman operator  $T$  is a **contraction mapping**
- 2 For any value functions  $u, v$ :  $\|Tu - Tv\|_\infty \leq \gamma \|u - v\|_\infty$
- 3 By Banach Fixed Point Theorem,  $T$  has unique fixed point  $v_*$
- 4 Value iteration:  $v_{k+1} = Tv_k$  converges to  $v_*$

## Rate of Convergence

$$\|v_k - v_*\|_\infty \leq \gamma^k \|v_0 - v_*\|_\infty \quad (8)$$

Geometric convergence with rate  $\gamma$

# Error Bounds and Stopping Criteria

**Practical Question:** When should we stop the algorithm?

## Error Bound

If  $\|v_{k+1} - v_k\|_\infty < \theta$ , then:

$$\|v_k - v_*\|_\infty \leq \frac{\theta}{1 - \gamma} \quad (9)$$

## Policy Loss Bound:

- Let  $\pi_k$  be greedy policy w.r.t.  $v_k$
- If  $\|v_k - v_*\|_\infty \leq \epsilon$ , then:

$$\|v_* - v_{\pi_k}\|_\infty \leq \frac{2\gamma\epsilon}{1 - \gamma} \quad (10)$$

## Practical Insight

Small value function errors lead to near-optimal policies



# Synchronous vs Asynchronous Updates

## Synchronous Value Iteration:

- Update all states simultaneously
- Use  $V_k$  to compute  $V_{k+1}$
- Requires two arrays
- Standard algorithm

### Pseudocode

```
for s in states:  
    V_new[s] = max_a Q(s,a)  
V = V_new
```

## Asynchronous Value Iteration:

- Update states one at a time
- Use most recent values
- In-place updates
- Often faster convergence

### Pseudocode

```
for s in states:  
    V[s] = max_a Q(s,a)  
// immediately use new V[s]
```

## Note

Both versions converge to  $v_*$ , but asynchronous often faster in practice

# Value Iteration Variants

## 1. Gauss-Seidel Value Iteration

- Update states in systematic order
- Use newest available values
- Faster convergence than Jacobi method

## 2. Prioritized Sweeping

- Update states with largest Bellman error first
- Maintain priority queue of states
- Much faster for sparse problems


## 3. Real-Time Value Iteration

- Update only visited states
- Useful for large state spaces
- Agent-centric approach

### Key Insight

All variants maintain convergence guarantees while improving efficiency

# Grid World Example

S1	S2	T
S3	X	S4
 S	S6	S7

## Problem Setup:

- $3 \times 3$  grid world
- Start: S5, Goal: T (reward +10)
- Obstacle: X (impassable)
- Actions: Up, Down, Left, Right
- Step reward: -1
- $\gamma = 0.9$

## Transition Model:

- 0.8 prob: intended direction
- 0.1 prob: each perpendicular direction
- Hit wall: stay in place

# Grid World: Value Iteration Steps

**Iteration 0:**

0	0	0
0	X	0
0	0	0

**Iteration 1:**

-1	-1	9
-1	X	9
-1	-1	-1

**Iteration 2:**

-1.9	6.4	9
-1.9	X	6.4
-1.9	-1.9	-1.9

**Iteration 5:**

4.6	6.4	9
2.3	X	6.4
0.4	2.3	4.6

## Value Update Example (S6):

$$V_1(S6) = \max\{Q(S6, \text{up}), Q(S6, \text{down}), Q(S6, \text{left}), Q(S6, \text{right})\} \quad (11)$$

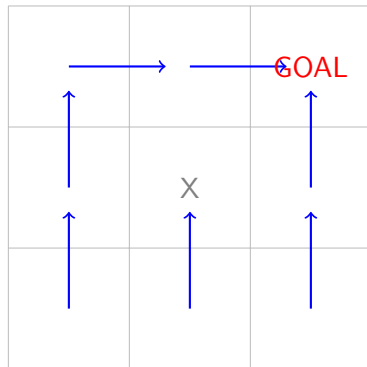
$$= \max\{-1, -1, -1, 6.4\} = 6.4 \quad (12)$$

Where:  $Q(S6, \text{right}) = -1 + 0.9 \times [0.8 \times 9 + 0.1 \times (-1) + 0.1 \times (-1)] = 6.4$

# Grid World: Optimal Policy

6.1	7.4	10
4.3	X	7.4
2.4	4.3	6.1

**Final Values:**



**Optimal Policy:**

**Policy Extraction:**

$$\pi_*(s) = \arg \max_a \sum_{s'} P(s'|s, a) [R + \gamma V(s')] \quad (13)$$

# Complexity Analysis

## Time Complexity per Iteration:

- For each state  $s$ :  $O(|A| \times |S|)$  operations
- Total per iteration:  $O(|S| \times |A| \times |S|) = O(|S|^2|A|)$

## Number of Iterations:

- Convergence rate:  $O(\log(1/\epsilon)/\log(1/\gamma))$
- For  $\epsilon = 10^{-6}$ ,  $\gamma = 0.9$ :  $\approx 131$  iterations

## Total Complexity:

$$O\left(\frac{|S|^2|A|}{\log(1/\gamma)} \log\left(\frac{1}{\epsilon}\right)\right) \quad (14)$$

## Curse of Dimensionality

For  $n$  state variables with  $k$  values each:  $|S| = k^n$  (exponential!)

# Memory Requirements

## Space Complexity:

- Value function:  $O(|S|)$
- Transition model:  $O(|S|^2|A|)$
- Policy:  $O(|S|)$

## Optimization Techniques:

- 1 **In-place updates:** Reduce memory by half
- 2 **Sparse representations:** For structured problems
- 3 **Function approximation:** For large state spaces
- 4 **State abstraction:** Group similar states

## Practical Limits

- Tabular VI:  $\sim 10^6$  states (modern computers)
- Beyond this: Need approximation methods

# Value Iteration: Strengths

## Advantages

- **Guaranteed convergence** to optimal solution
- **Simple to implement** and understand
- **No policy needed** during learning
- **Anytime algorithm**: can stop early for approximate solution
- **Parallelizable**: states can be updated independently
- **Memory efficient**: only needs value function

## Mathematical Elegance

- Direct implementation of Bellman optimality equation
- Clear convergence theory and error bounds
- Foundation for many advanced RL algorithms



# Value Iteration: Limitations

## Major Limitations

- **Requires complete MDP model** ( $P, R$ )
- **Curse of dimensionality:**  $O(|S|^2|A|)$  per iteration
- **Only works for finite state/action spaces**
- **Slow for large discount factors**  $\gamma \approx 1$
- **No online learning:** must know environment model

## When VI Struggles:

- Continuous state/action spaces
- Unknown environment dynamics
- Very large state spaces ( $> 10^6$  states)
- Real-time applications requiring fast responses

## Solutions

Function approximation, sampling methods, model-free approaches

# Value Iteration vs Policy Iteration

Aspect	Value Iteration	Policy Iteration
Updates	Value function only	Policy + Value function
Convergence	$O(\log(1/\epsilon))$	Fewer iterations
Per iteration	$O( S ^2 A )$	$O( S ^3 +  S ^2 A )$
Memory	$O( S )$	$O( S )$
Implementation	Simpler	More complex
Early stopping	Good approx. policy	Poor policy

## When to use each:

- **Value Iteration:** When you want simplicity, anytime behavior
- **Policy Iteration:** When exact solution needed, small state spaces

# Value Iteration Extensions

## 1. Approximate Value Iteration

- Use function approximation:  $V(s) \approx \hat{V}(s; \theta)$
- Neural networks, linear functions, etc.
- Enables large/continuous state spaces

## 2. Asynchronous Value Iteration

- Update states in any order
- Real-time dynamic programming
- Prioritized sweeping variants

## 3. Multi-objective Value Iteration

- Vector-valued rewards and values
- Pareto-optimal policies
- Applications in robotics, finance

# Real-World Applications

## 1. Game Playing

- Chess, Go, Backgammon endgames
- Perfect information games
- Tablebase generation

## 2. Robotics

- Path planning in discrete grids
- Navigation with known maps
- Manipulation planning

## 3. Operations Research

- Inventory management
- Queueing systems optimization
- Maintenance scheduling

## 4. Finance

- Portfolio optimization
- Options pricing (American options)
- Risk management

# Key Takeaways

## Value Iteration Algorithm

- **Core idea:** Iteratively apply Bellman optimality operator
- **Convergence:** Guaranteed for finite MDPs with  $\gamma < 1$
- **Complexity:**  $O(|S|^2|A|)$  per iteration

## Theoretical Foundation

- Based on contraction mapping theorem
- Geometric convergence rate  $\gamma$
- Clear error bounds and stopping criteria

## Practical Considerations

- Works well for moderate-sized problems
- Requires complete environment model
- Foundation for more advanced RL methods

## What's Coming Next:

- **Policy Iteration:** Alternative DP approach
- **Generalized Policy Iteration:** Unified framework
- **Monte Carlo Methods:** Model-free approaches
- **Temporal Difference Learning:** Online model-free methods

## Practice Problems:

- Implement value iteration for small grid worlds
- Analyze convergence for different  $\gamma$  values
- Compare with policy iteration on same problems
- Explore asynchronous update schemes

 Questions?