

INTRO TO RL

Xudong Wang

¹Department of Industrial & Systems Engineering, The University of Tennessee, Knoxville.



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

DEFINITION

- The State Value Function, denoted as $V(s)$, is a fundamental concept in reinforcement learning. It represents the expected cumulative reward an agent can obtain from a given state s while following a specific policy.
- Mathematical Representation:

$$V^{\pi}(s) = E_{\pi}[G_t | S_t = s] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right], \forall s \in \mathcal{S}$$

UNDERSTANDING STATE VALUE FUNCTION

- $V(s)$ quantifies the goodness of a state s under a given policy π .
- It indicates how desirable it is to be in state s when following policy π .

State Value Function is used for:

- Evaluating the quality of a state **under a policy**.
- Solving reinforcement learning problems using methods like Value Iteration and Policy Iteration.

PROPERTIES OF STATE VALUE FUNCTION

- **Markov Property:** $V(s)$ depends only on the current state s and the policy π .
- **Recursive Nature:** $V(s)$ is defined in terms of itself, involving future rewards.
- **Optimal State Values:** There exists an optimal policy π^* such that $V^*(s) = V^{\pi^*}(s)$, where $V^*(s)$ is the state value function for the optimal policy.

Computing State Value Function

- **Value Iteration:** An iterative method for finding the optimal $V^*(s)$.
- **Policy Evaluation:** Part of the Policy Iteration process.

EXAMPLE: GRID WORLD

- Consider a 5*5 grid world with states.
- An agent can move up, down, left, or right.
- Every move will have a -1 reward.
- When entering location A, it will have a +10 reward (no move cost).
- When entering location B, it will have a -2 reward (no move cost). A is the terminal point.
- The goal is to reach the terminal state with maximum cumulative reward.

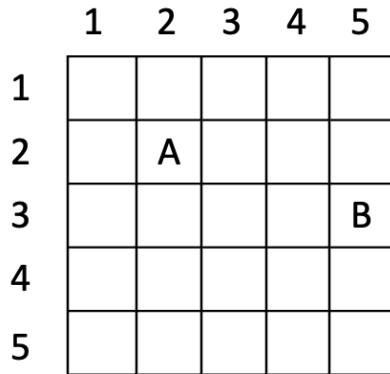


FIGURE: Grid World Example

COMPUTING STATE VALUE FUNCTION

	1	2	3	4	5
1	0	0	0	0	0
2	0	T	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0
5	0	0	0	0	0

→

	1	2	3	4	5
1	-1	10	-1	-1	-1
2	10	T	10	-1	-1
3	-1	10	-1	-1	-1
4	-1	-1	-1	-1	-1
5	-1	-1	-1	-1	-1

$$\begin{aligned}
 \pi: \text{greedy} \quad V(1, 2) &= \max(E(\text{left}), E(\text{right}), E(\text{down})) \\
 \gamma: 0.9 \quad &= \max(-1 + \gamma V(1, 1), -1 + \gamma V(3, 1), 10) \\
 &= 10
 \end{aligned}$$

$$\begin{aligned}
 V(3, 4) &= \max(E(\text{left}), E(\text{right}), E(\text{down}), E(\text{up})) \\
 &= \max(-1 + \gamma V(3, 3), -2 + \gamma V(3, 5), -1 + \gamma V(4, 4), -1 + \gamma V(2, 4)) \\
 &= -1
 \end{aligned}$$

FIGURE: State Value Function for Grid World

COMPUTING STATE VALUE FUNCTION

	1	2	3	4	5
1	0	0	0	0	0
2	0	T	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0
5	0	0	0	0	0

→

	1	2	3	4	5
1	-1	2.67	-1	-1	-1
2	2.67	T	2.33	-1	-1.34
3	-1	2.33	-1	-1.25	-1
4	-1	-1	-1	-1	-1.34
5	-1	-1	-1	-1	-1

$$\begin{aligned}
 \pi: \text{random} \quad \gamma: 0.9 \quad V(1,2) &= \frac{1}{3} * (E(\text{left}) + E(\text{right}) + E(\text{down})) \\
 &= \frac{1}{3} * (-1 + \gamma V(1,1) - 1 + \gamma V(3,1) + 10) \\
 &= \frac{8}{3}
 \end{aligned}$$

$$\begin{aligned}
 V(3,4) &= \frac{1}{4} * (E(\text{left}) + E(\text{right}) + E(\text{down}) + E(\text{up})) \\
 &= \frac{1}{4} * (-1 + \gamma V(3,3) - 2 + \gamma V(3,5) - 1 + \gamma V(4,4) - 1 + \gamma V(2,4)) \\
 &= -1.25
 \end{aligned}$$

ACTION VALUE FUNCTION (Q-FUNCTION)

- The Action Value Function, denoted as $Q(s, a)$, represents the expected cumulative reward an agent can obtain from taking action a in state s and following a specific policy.
- Mathematical Representation:

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

UNDERSTANDING ACTION VALUE FUNCTION

- $Q(s, a)$ quantifies the goodness of taking a specific action a in a given state s under a particular policy π .
- It indicates the expected return if the agent selects action a while in state s and continues to follow policy π .

Action Value Function is used for:

- Evaluating the quality of actions in specific states.
- Selecting the best action to maximize the expected return.

RELATIONSHIP BETWEEN STATE AND ACTION VALUE FUNCTIONS

Connection:

The State Value Function $V(s)$ based on greedy method and Action Value Function $Q(s, a)$ are related as follows:

$$V(s) = \max_a Q(s, a)$$

Optimal Policy: For the optimal policy π^* :

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

where $Q^*(s, a)$ is the action value function for the optimal policy.

Similar to State Value Function, Action Value Function can be computed using methods like Value Iteration and Policy Iteration.

COMPUTING ACTION VALUE FUNCTION

Based on a random policy:

$$\begin{aligned} Q^\pi((0,0), right) &= \sum_{s'} p(s', r|s, a)[r + \gamma \sum_{a'} \pi(a'|s') Q^\pi(s', a')] \\ &= \sum_{s'} p(s', r|s, a)[r + \gamma V^\pi(s')] \end{aligned}$$

$$s' = (0,0), r = -1, \gamma = 0.9$$

COMPUTING ACTION VALUE FUNCTION

Consider another example: The agent is in state 1. When it takes action a , it will get a reward α and reach state 2 or 3 with an equal probability. When it takes action b , it will have a 0.3 probability of getting a reward β and reaching state 3; and a 0.7 probability of getting a reward v and reaching state 4.

$$\begin{aligned} Q^\pi(1, a) &= \sum_{s'} p(s', r|1, a)[r + \gamma V^\pi(s')] \\ &= 0.5 * [\alpha + \gamma V^\pi(2)] + 0.5 * [\alpha + \gamma V^\pi(3)] \end{aligned}$$

$$\begin{aligned} Q^\pi(1, b) &= \sum_{s'} p(s', r|1, a)[r + \gamma V^\pi(s')] \\ &= 0.3 * [\beta + \gamma V^\pi(3)] + 0.7 * [v + \gamma V^\pi(4)] \end{aligned}$$

DYNAMIC PROGRAMMING

Dynamic programming is the collection of algorithms that can be used to compute optimal policies given perfect model of the environment.

- DP constitutes a theoretically optimal methodology
- In reality it is often limited since DP is computationally expensive.
- Key Idea (In general): Use value function to derive optimal policy
- Need to find a way to calculate optimal state value function

VALUE ITERATION

Definition:

Value Iteration is an iterative algorithm used to find the optimal state value function $V^*(s)$ and, consequently, the optimal policy π^* in reinforcement learning.

- It starts with an initial value function and repeatedly updates it until it converges to the optimal value function.
- The optimal policy can be derived directly from the optimal value function.

The Value Iteration algorithm works as follows:

- Initialize $V(s)$ arbitrarily for all states.
- Iterate until convergence:
 - Update each state's value using the Bellman optimality equation:

$$V(s) \leftarrow \max_a \sum_{s',r} P(s',r|s,a)[r + \gamma V(s')]$$

- The resulting $V^*(s)$ represents the optimal state values.

VALUE ITERATION

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

```
|  $\Delta \leftarrow 0$   
| Loop for each  $s \in \mathcal{S}$ :  
|    $v \leftarrow V(s)$   
|    $V(s) \leftarrow \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$   
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
until  $\Delta < \theta$ 
```

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$$

POLICY ITERATION

Definition Policy Iteration is an algorithm used to find the optimal policy in reinforcement learning. It alternates between two steps: policy evaluation and policy improvement.

- In the policy evaluation step, it evaluates the current policy to get its value function $V^\pi(s)$.
- In the policy improvement step, it improves the policy based on the value function.
- It repeats these steps until the policy converges to the optimal policy.

The Policy Iteration algorithm works as follows:

- Initialize a policy π arbitrarily.
- Iterate until convergence:
 - Policy Evaluation: Compute the value function $V^\pi(s)$ for the current policy.
 - Policy Improvement: Update the policy based on the current value function:

$$\pi(s) \leftarrow \arg \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V^\pi(s')]$$

- The resulting policy π^* is the optimal policy.

POLICY ITERATION

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow *true*

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow *false*

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

VALUE ITERATION VS. POLICY ITERATION

Comparison

■ Value Iteration:

- More straightforward and computationally efficient.
- Directly computes the optimal state values and policy.
- May require more iterations to converge.

■ Policy Iteration:

- Faster convergence in practice.
- Alternates between policy evaluation and policy improvement.
- Requires solving a linear system of equations in the policy evaluation step.

Choice of Method The choice between Value Iteration and Policy Iteration often depends on the specific problem, computational resources, and the need for a fast or efficient solution.

COMPARISON BETWEEN DP AND MC

DP

- theoretically optimal result
- expensive computation
- model-based

MC

- can be “just as well” as DP
- less computation and memory
- model free

MONTE CARLO POLICY EVALUATION

- We still need to estimate $V(s)$. Estimate from experience
 - Average all returns observed after visiting a given state.
 - Each occurrence of state s in an episode is called a visit to s
- Generate an episode following π :
 $s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T$
- First visit MC: Method averages just the returns following first visits to s
- Every Visit MC: Method averages the returns following all the visits to s

MONTE CARLO POLICY EVALUATION

First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy π to be evaluated

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

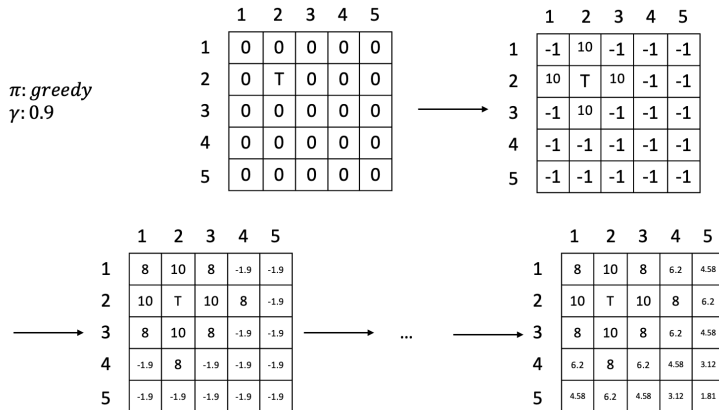
Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

Append G to $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

EXAMPLE OF VALUE ITERATION

Converge after 6 iterations:



EXAMPLE OF POLICY ITERATION

Converge after 4 iterations:

