

Exact Methods : Value Iteration

Easwar Subramanian

TCS Innovation Labs, Hyderabad

Email : easwar.subramanian@tcs.com / cs5500.2020@iith.ac.in

August 26, 2022

- 1 Review
- 2 Value Iteration
- 3 Possible Extensions
- 4 Few Remarks

Review

Question : Is there a way to arrive at π_* starting from an arbitrary policy π ?

Answer : Policy Iteration

► **Evaluate** the policy π

★ Compute $V^\pi(s) = \mathbb{E}_\pi(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s)$

► **Improve** the policy π

$$\pi'(s) = \text{greedy}(V^\pi(s))$$

$$\pi_0 \xrightarrow{\text{E}} V^{\pi_0} \xrightarrow{\text{I}} \pi_1 \xrightarrow{\text{E}} V^{\pi_1} \xrightarrow{\text{I}} \pi_2 \xrightarrow{\text{E}} \dots \xrightarrow{\text{I}} \pi^* \xrightarrow{\text{E}} V^*,$$

Algorithm Policy Iteration

- 1: Start with an initial policy π_1
- 2: **for** $i = 1, 2, \dots, N$ **do**
- 3: Evaluate $V^{\pi_i}(s) \quad \forall s \in \mathcal{S}$. That is,
- 4: **for** $k = 1, 2, \dots, K$ **do**
- 5: For all $s \in \mathcal{S}$ calculate

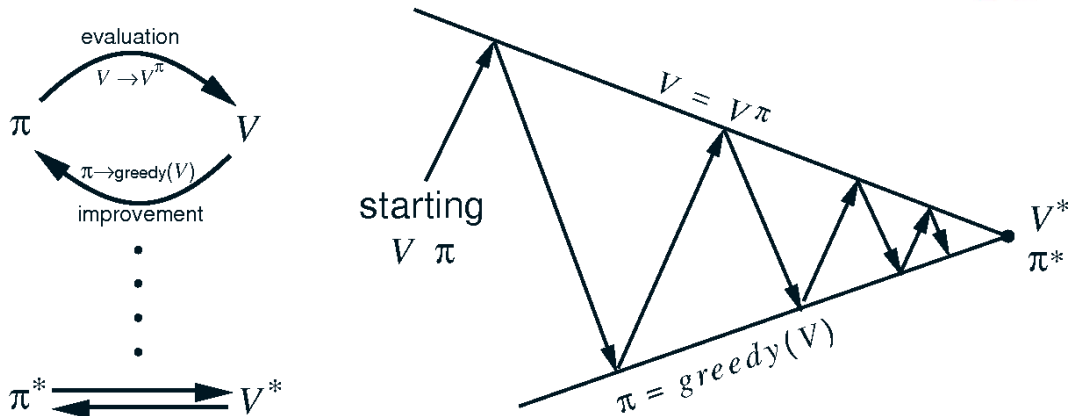
$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k^{\pi_i}(s')]$$

- 6: **end for**
- 7: Perform policy Improvement

$$\pi_{i+1} = \text{greedy}(V^{\pi_i})$$

- 8: **end for**
-

Policy Iteration : Schematic Representation



- ▶ The sequence $\{\pi_1, \pi_2, \dots\}$ is guaranteed to converge.
- ▶ At convergence, both current policy and the value function associated with the policy are optimal.

Value Iteration

Question : Is there a way to arrive at V_* starting from an arbitrary value function V_0 ?

Answer : Value Iteration

Recall the Bellman Evaluation Equation for an MDP with policy π

$$V^\pi(s) = \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

Question : Can we have a recursive formulation for $V_*(s)$?

$$V_*(s) = \max_a Q_*(s, a) = \max_a \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma V_*(s')) \right]$$

Similarly, there is a recursive formulation for $Q_*(\cdot, \cdot)$

$$Q_*(s, a) = \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \left(\mathcal{R}_{ss'}^a + \gamma \max_{a'} Q_*(s', a') \right) \right]$$

Question : These are also a system of equations with $n = |\mathcal{S}|$ with n variables. Can we solve them ?

Answer : Optimality equations are **non-linear** system of equations with n unknowns and n non-linear constraints (i.e., the max operator).

Solving the Bellman Optimality Equation

- ▶ Bellman optimality equations are non-linear
- ▶ In general, there are no closed form solutions
- ▶ Iterative methods are typically used

Principle of Optimality

The tail of an optimal policy must be optimal

- Any optimal policy can be subdivided into two components; an optimal first action, followed by an optimal policy from successor state s' .

Bellman optimality equation :

$$V_*(s) = \max_a \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma V_*(s')) \right]$$

Optimal Substructure : Optimal solution can be constructed from optimal solutions to subproblems

Overlapping Subproblems : Problem can be broken down into subproblems and can be reused several times

- ▶ Markov Decision Processes, generally, satisfy both these characteristics
- ▶ Dynamic Programming is a popular solution method for problems having such properties

- ▶ Suppose we know the value $V_*(s')$
- ▶ Then the solution $V_*(s)$ can be found by one step look ahead

$$V_*(s) \leftarrow \max_a \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma V_*(s')) \right]$$

- ▶ Idea of value iteration is to perform the above updates iteratively

Algorithm Value Iteration

1: Start with an initial value function $V_0(\cdot)$;

2: **for** $k = 0, 1, 2, \dots, K$ **do**

3: **for** $s \in \mathcal{S}$ **do**

4: Calculate

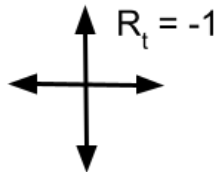
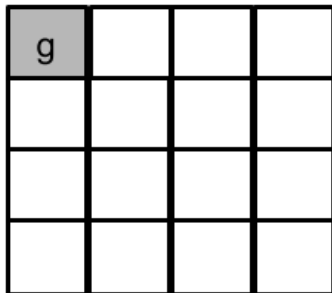
$$V_{k+1}(s) \leftarrow \max_a \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma V_k(s')) \right]$$

5: **end for**

6: **end for**

Value Iteration : Example

No noise and discount factor $\gamma = 1$



Value Iteration : Example

$$V_{k+1}(s) \leftarrow \max_a \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a (\mathcal{R}_{ss'}^a + \gamma V_k(s')) \right]$$

g			

Problem

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

V_1

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1

V_2

0	-1	-2	-2
-1	-2	-2	-2
-2	-2	-2	-2
-2	-2	-2	-2

V_3

0	-1	-2	-3
-1	-2	-3	-3
-2	-3	-3	-3
-3	-3	-3	-3

V_4

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-4
-3	-4	-4	-4

V_5

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-5

V_6

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-6

V_7

- ▶ The sequence of value functions $\{V_1, V_2, \dots\}$ converge
- ▶ It converges to V_*
- ▶ Convergence is independent of the choice of V_0 .
- ▶ Intermediate value functions need not correspond to a policy in the sense of satisfying the Bellman Evaluation Equation
- ▶ However, for any k , one can come up with a greedy policy as follows

$$\pi_{k+1}(s) \leftarrow \text{greedy} V_k(s)$$

There is a recursive formulation for $Q_*(\cdot, \cdot)$

$$Q_*(s, a) = \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \left(\mathcal{R}_{ss'}^a + \gamma \max_{a'} Q_*(s', a') \right) \right]$$

One could similarly conceive an iterative algorithm to compute optimal Q_* using the above recursive formulation !!

Iterative application of Bellman Evaluation Equation

Iterative update rule :

$$V_{k+1}^{\pi}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k^{\pi}(s')]$$

The sequence of value functions $\{V_1^{\pi}, V_2^{\pi}, \dots, \}$ converge to V^{π}

Policy Iteration : Example Revisited

Update Rule :

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k^{\pi_i}(s')]$$

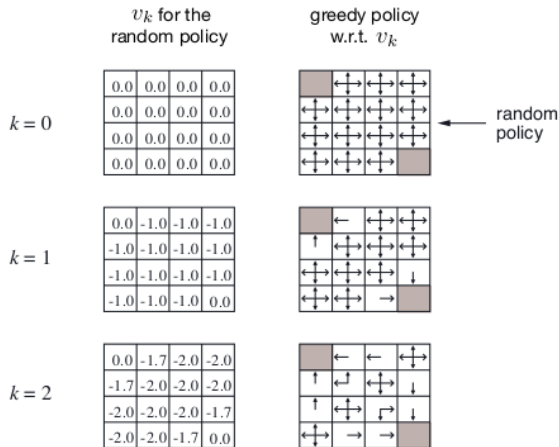
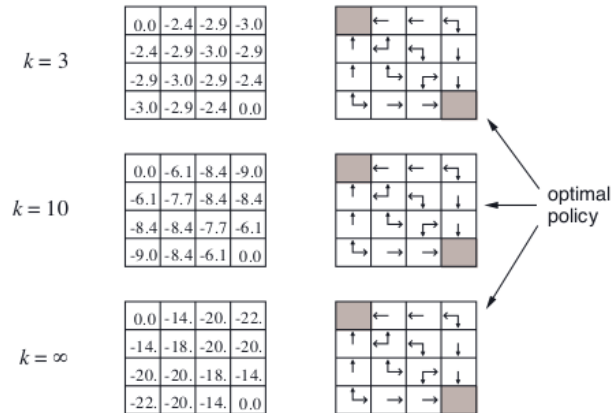


Figure Source: David Silver's UCL course

Policy Iteration : Example Revisited



Can we computationally simplify policy iteration process ?

- ▶ We need not wait for policy evaluation to converge to V^π
- ▶ We can have a stopping criterion like ϵ -convergence of value function evaluation or K iterations of policy evaluation
- ▶ Extreme case of $K = 1$ is **value iteration**. We update the policy every iteration

- ▶ Dynamic Programming assumes full knowledge of MDP
- ▶ Used for both **prediction** and **control** in an MDP
- ▶ Prediction
 - ★ Input MDP ($\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$) and policy π
 - ★ Output : $V^\pi(\cdot)$
- ▶ Control
 - ★ Input MDP ($\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$)
 - ★ Output : Optimal value function $V_*(\cdot)$ or optimal policy π_*

Problem	Bellman Equation	Algorithm
Prediction	Bellman Evaluation Equation	Policy Evaluation
Control	Bellman Evaluation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

- ▶ All the methods described above have synchronous backups
- ▶ All states are backed up in every iteration

Possible Extensions

- ▶ Updates to states are done individually, in any order
- ▶ For each selected state, apply the appropriate backup
- ▶ Can significantly reduce computation
- ▶ Convergence guarantees exist, if all states are selected sufficient number of times

- ▶ Idea : update only states that are relevant to agent
- ▶ After each time step, we get s_t, a_t, r_{t+1}
- ▶ Perform the following update

$$V(s_t) \leftarrow \max_a \left[\sum_{s' \in \mathcal{S}} \mathcal{P}_{s_t s'}^a (\mathcal{R}_{s_t s'}^a + \gamma V(s')) \right]$$

Few Remarks

- ▶ **MDP Setting** : The agent has knowledge of the state transition matrices $\mathcal{P}_{ss'}^a$ and the reward function \mathcal{R}
- ▶ **RL Setting** : The agent does not have knowledge of the state transition matrices $\mathcal{P}_{ss'}^a$ and the reward function \mathcal{R}
 - ★ The goal in both cases are same; Determine optimal sequence of actions such that the total discounted future reward is maximum.
 - ★ Although, this course would assume Markovian structure to state transitions, in many (sequential) decision making problems we may have to consider the history as well.

- ▶ Recall that a (stochastic) policy is a distribution over actions given states
- ▶ Markov policy means that the policy depends only on the current state and not on the history
- ▶ Policies could be stationary or non-stationary
- ▶ In general, the optimal policy for an MDP need not be unique
- ▶ For finite horizon MDP, the optimal policy need not be even stationary
- ▶ For infinite horizon, an MDP admits an optimal policy that is deterministic and stationary. But there could other optimal policies that are stochastic and non-stationary.

- ▶ The grid world problem is an example **stochastic shortest path** problem where we consider only policies that are 'proper'

- ★ A policy that has a non-zero chance to finally reach the terminal state

Under this assumption the theory on convergence will work out for even $\gamma = 1$.

- ▶ The total discounted return G_t could have infinite terms or $\gamma = 1$ but not both

- ▶ How do we know that value iteration converges to V_* ?
- ▶ Or that iterative policy evaluation converges to V_π ?
- ▶ And therefore that policy iteration converges to π_* ?
- ▶ Is the solution unique ?
- ▶ How fast do these algorithms converge ? (Depends on discount factor γ)
- ▶ These questions were resolved by
 - ★ **Banach Fixed Point Theorem / Contraction Mapping Theorem**