

AIL 722: Reinforcement Learning

Lecture 10: Policy Iteration and Value Iteration

Raunak Bhattacharyya



ScAI

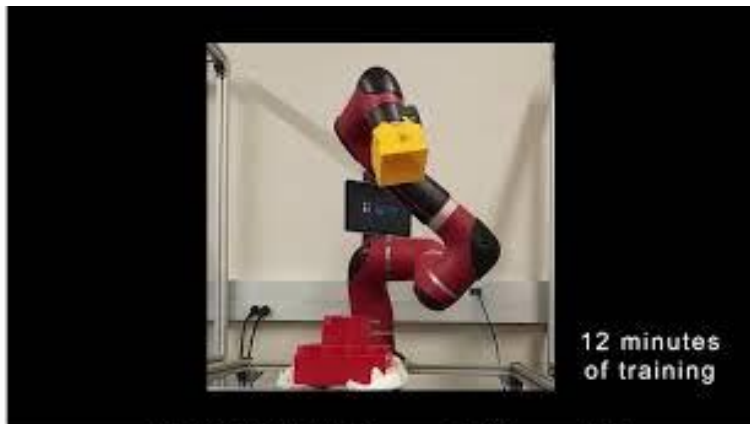
YARDI SCHOOL OF ARTIFICIAL INTELLIGENCE
INDIAN INSTITUTE OF TECHNOLOGY DELHI

Outline

- Defining Optimality
- Policy Iteration
- Value Iteration

Discounting

Episodic



© Authors of ICRA 2018 Paper 1799

Thu AM

Pod Q.2

[Source: Youtube](#)

Infinite horizon



[Source: Youtube](#)

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{p_{\theta}(\tau)} \left[\sum_{t=0}^T \gamma^t \cdot r(s_t, a_t) \right]$$

Bellman Equation

$$V^\pi(s) = r(s, \pi(s)) + \gamma \cdot \mathbb{E}_{p(s'|s, \pi(s))} \left[V^\pi(s') \right]$$

Discounting

Deterministic policies

Optimality

- Goal: Finding a policy that achieves a lot of reward over the long run
- Notion of betterness: A policy is better than or equal to another policy if its expected return is greater than or equal to that of the other policy for all states

$$\pi \geq \pi' \text{ if and only if } V^\pi(s) \geq V^{\pi'}(s) \quad \forall s \in \mathcal{S}$$

- There is always at least one policy that is better than or equal to all the other policies. This is called an optimal policy

Optimality

- All the optimal policies share the same state value function as well as the same state-action value function

$$V^*(s) = \max_{\pi} V^{\pi}(s) \quad \forall s \in \mathcal{S}$$

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) \quad \forall s \in \mathcal{S} \text{ and } \forall a \in \mathcal{A}$$

$$J(\theta) = \mathbb{E}_{p(s_1)} \left[V^{\pi}(s_1) \right]$$

Policy Iteration



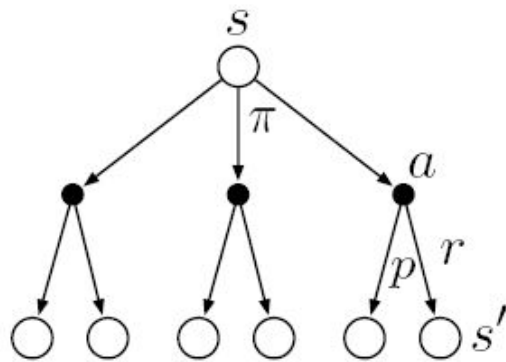
1. Evaluate $V^\pi(s)$

2. Set $\pi \leftarrow \pi_{\text{new}}$

$$\pi_{\text{new}} = \begin{cases} 1 & \text{if } a = \arg \max_a A^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Policy Evaluation

$$V^\pi(s) = r(s, \pi(s)) + \gamma \cdot \mathbb{E}_{p(s'|s, \pi(s))} \left[V^\pi(s') \right]$$



Backup diagram for v_π

Iterative Policy Evaluation

Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input π , the policy to be evaluated

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 $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

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

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

until $\Delta < \theta$

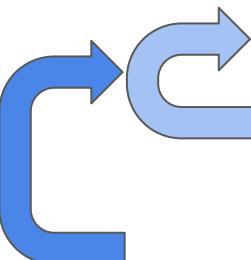
Demo: Iterative Policy Evaluation

GridWorld: Dynamic Programming Demo

Policy Evaluation (one sweep) Policy Update Toggle Value Iteration Reset

0.00 ↖	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↗
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕					0.00 ↕				0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0		0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕ R -1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕ R 1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕ R -1.0	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0		0.00 ↕ R -1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕

Policy Iteration



1. $V^\pi(s) = r(s, \pi(s)) + \gamma \cdot \mathbb{E}_{p(s'|s, \pi(s))} [V^\pi(s')]$
2. Set $\pi \leftarrow \pi_{\text{new}}$

$$\pi_{\text{new}} = \begin{cases} 1 & \text{if } a = \arg \max_a A^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Finding the Policy

$$\pi_{\text{new}} = \begin{cases} 1 & \text{if } a = \arg \max_a A^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$$

$$A^\pi(s, a) = r(s, a) + \mathbb{E}_{p(s'|s, a)} \left[V^\pi(s') \right] - V^\pi(s)$$

Goal: Find the argmax

Finding the Policy

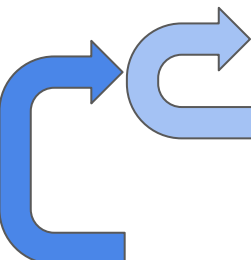
Goal: Find the argmax

$$\pi_{\text{new}} = \begin{cases} 1 & \text{if } a = \arg \max_a Q^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$$

$$Q^\pi(s, a) = r(s, a) + \gamma \cdot \mathbb{E}_{p(s'|s,a)} \left[V^\pi(s') \right]$$

Extract the policy using Q-function (table in case of tabular setting)

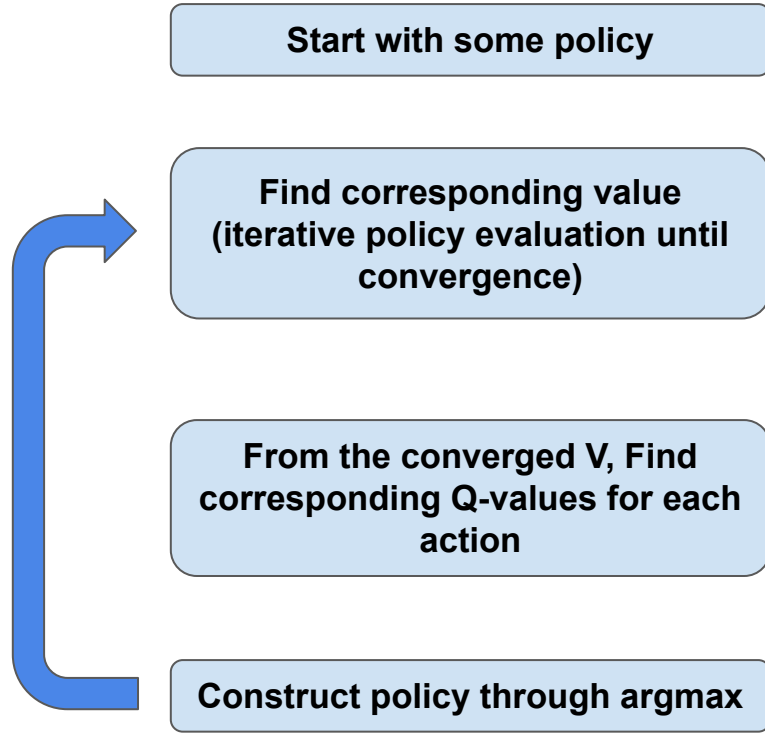
Policy Iteration: Using Q-values



1. $V^\pi(s) = r(s, \pi(s)) + \gamma \cdot \mathbb{E}_{p(s'|s, \pi(s))} [V^\pi(s')]$
2. Set $\pi \leftarrow \pi_{\text{new}}$

$$\pi_{\text{new}} = \begin{cases} 1 & \text{if } a = \arg \max_a Q^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Workflow



Pseudocode?

Policy Iteration Demo

GridWorld: Dynamic Programming Demo

Policy Evaluation (one sweep) Policy Update Toggle Value Iteration Reset

0.00 ↖	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↗
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕					0.00 ↕				0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0		0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕ R -1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕ R 1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕ R -1.0	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0		0.00 ↕ R -1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕

Value Iteration

$a_?$
$a_?$
$a_?$
$a_?$
$a_?$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

a_2
a_1
a_3
a_3
a_1

$a_?$
$a_?$
$a_?$
$a_?$
$a_?$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

a_2
a_1
a_3
a_3
a_1


$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

Value Iteration

Start with a random value function $V(s)$

- 
1. Set $Q(s, a) \leftarrow r(s, a) + \gamma \cdot \mathbb{E}_{p(s'|s,a)} \left[V^\pi(s') \right]$
 2. Set $V(s) \leftarrow \max_a Q(s, a)$

Pseudocode?

How do we find the policy?

Value Iteration Demo

GridWorld: Dynamic Programming Demo

Policy Evaluation (one sweep) Policy Update Toggle Value Iteration Reset

0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖
0.00 ↖					0.00 ↖				0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖ R -1.0		0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖		0.00 ↖ R -1.0	0.00 ↖ R -1.0	0.00 ↖	0.00 ↖	0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖		0.00 ↖ R 1.0	0.00 ↖ R -1.0	0.00 ↖	0.00 ↖ R -1.0	0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖		0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖ R -1.0	0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖ R -1.0		0.00 ↖ R -1.0	0.00 ↖ R -1.0	0.00 ↖	0.00 ↖	0.00 ↖
0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖	0.00 ↖