

Reinforcement Learning

Lecture 2b:

Value Iteration

[SutBar] Sec. 4.1, 4.1, [Sze] Sec. 2.2, 2.3,
[Put] Sec. 6.1-6.3, [SigBuf] Chap. 1

Outline

- Convergence properties of
 - Policy evaluation
 - Value iteration

Value Iteration Algorithm

valueiteration(MDP)

$$V_0^*(s) \leftarrow \max_a R(s, a) \quad \forall s$$

For $t = 1$ to h do

$$V_t^*(s) \leftarrow \max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V_{t-1}^*(s') \quad \forall s$$

Return V^*

Optimal policy π^*

$$t = 0: \pi_0^*(s) \leftarrow \arg\max_a R(s, a) \quad \forall s$$

$$t > 0: \pi_t^*(s) \leftarrow \arg\max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V_{t-1}^*(s') \quad \forall s$$

NB: t indicates the # of time steps to go (till end of process)

π^* is **non stationary** (i.e., time dependent)

Value Iteration

- Matrix form:

R^a : $|S| \times 1$ column vector of rewards for a

V_t^* : $|S| \times 1$ column vector of state values

T_a^t : $|S| \times |S|$ matrix of transition prob. for a

valueiteration(MDP)

$V_0^* \leftarrow \max_a R^a$

For $t = 1$ to h do

$V_t^* \leftarrow \max_a R^a + \gamma T^a V_{t-1}^*$

Return V^*

Example: MDP $|A|=|S|=2$

2 states \therefore 2-action MDP

$$T^{a_1} = \begin{matrix} & \begin{matrix} s_1 & s_2 \end{matrix} \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 0.3 & 0.7 \\ 0.8 & 0.2 \end{pmatrix} \end{matrix}$$

$$T^{a_2} = \begin{matrix} & \begin{matrix} s_1 & s_2 \end{matrix} \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 0.7 & 0.3 \\ 0.2 & 0.8 \end{pmatrix} \end{matrix}$$

$$R^{a_1} = \begin{matrix} s_1 & 0 \\ s_2 & 10 \end{matrix} \quad R^{a_2} = \begin{matrix} s_1 & -5 \\ s_2 & +5 \end{matrix}$$

$$R^{a_1} + \gamma T^{a_1} V^*$$

$$\max = \left\{ \begin{pmatrix} 0 \\ 10 \end{pmatrix} + 0.9 * \begin{pmatrix} 0.3 & 0.7 \\ 0.8 & 0.2 \end{pmatrix} \begin{pmatrix} V^*(s_1) \\ V^*(s_2) \end{pmatrix}, \right.$$

$$\left. \begin{pmatrix} -5 \\ +5 \end{pmatrix} + 0.9 * \begin{pmatrix} 0.7 & 0.3 \\ 0.2 & 0.8 \end{pmatrix} \begin{pmatrix} V^*(s_1) \\ V^*(s_2) \end{pmatrix} \right\}$$

element - max .

Infinite Horizon

- Let $h \rightarrow \infty$
- Then $V_h^\pi \rightarrow V_\infty^\pi$ and $V_{h-1}^\pi \rightarrow V_\infty^\pi$

- **Policy evaluation:**

$$V_\infty^\pi(s) = R(s, \pi_\infty(s)) + \gamma \sum_{s'} \Pr(s'|s, \pi_\infty(s)) V_\infty^\pi(s') \quad \forall s$$

- **Bellman's equation:**

$$V_\infty^\pi(s) = \max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V_\infty^\pi(s')$$

Policy evaluation

- Linear system of equations

$$V_{\infty}^{\pi}(s) = R(s, \pi_{\infty}(s)) + \gamma \sum_{s'} \Pr(s'|s, \pi_{\infty}(s)) V_{\infty}^{\pi}(s') \quad \forall s$$

- Matrix form:

R : $|S| \times 1$ column vector of state rewards for π

V : $|S| \times 1$ column vector of state values for π

T : $|S| \times |S|$ matrix of transition prob for π

$$V = R + \gamma TV$$

Handwritten notes illustrating policy evaluation with a fixed policy:

$\pi(s_1) = a_1, \quad \pi(s_2) = a_2 \rightarrow \text{fixed policy}$

$R^{\pi} = \begin{pmatrix} 0 \\ 5 \end{pmatrix}$

$T^{\pi} = \begin{matrix} & \begin{matrix} s_1 & s_2 \end{matrix} \\ \begin{matrix} s_1 \\ s_2 \end{matrix} & \begin{pmatrix} 0.3 & 0.7 \\ 0.2 & 0.8 \end{pmatrix} \end{matrix}$

Annotations: $\leftarrow a_1$ (pointing to the first column) and $\rightarrow a_2$ (pointing to the second column).

Solving linear equations

- Linear system: $V = R + \gamma TV$
- Gaussian elimination: $(I - \gamma T)V = R$
- Compute inverse: $V = (I - \gamma T)^{-1}R$
- Iterative methods
 - Value iteration (a.k.a. Richardson iteration)
 - Repeat $V \leftarrow R + \gamma TV$

Contraction

- Let $H(V) \stackrel{\text{def}}{=} R + \gamma TV$ be the policy eval operator
- **Lemma 1:** H is a **contraction mapping**.

$$\|H(\tilde{V}) - H(V)\|_{\infty} \leq \gamma \|\tilde{V} - V\|_{\infty}$$

- **Proof** $\|H(\tilde{V}) - H(V)\|_{\infty}$
 - $= \|R + \gamma T\tilde{V} - R - \gamma TV\|_{\infty}$ (by definition)
 - $= \|\gamma T(\tilde{V} - V)\|_{\infty}$ (simplification)
 - $= \gamma \|T\|_{\infty} \|\tilde{V} - V\|_{\infty}$ (since $\|AB\| \leq \|A\| \|B\|$)
 - $= \gamma \|\tilde{V} - V\|_{\infty}$ (since $\max_s \sum_{s'} T(s, s') = 1$)

Convergence

- **Theorem 2:** Policy evaluation converges to V^π for any initial estimate V

$$\lim_{n \rightarrow \infty} H^{(n)}(V) = V^\pi \quad \forall V$$

- Proof
 - By definition $V^\pi = H^{(\infty)}(0)$, but policy evaluation computes $H^{(\infty)}(V)$ for any initial V
 - By Lemma 1, $\|H^{(n)}(V) - H^{(n)}(\tilde{V})\|_\infty \leq \gamma^n \|V - \tilde{V}\|_\infty$
 - Hence, when $n \rightarrow \infty$, then $\|H^{(n)}(\tilde{V}) - H^{(n)}(V)\|_\infty \rightarrow 0$ and $H^{(\infty)}(V) = V^\pi \quad \forall V$

WHAT IS THE DIFFERENCE BETWEEN A THEOREM, A LEMMA, AND A COROLLARY?

PROF. DAVE RICHESON

- (1) **Definition**—a precise and unambiguous description of the meaning of a mathematical term. It characterizes the meaning of a word by giving all the properties and only those properties that must be true.
- (2) **Theorem**—a mathematical statement that is proved using rigorous mathematical reasoning. In a mathematical paper, the term theorem is often reserved for the most important results.
- (3) **Lemma**—a minor result whose sole purpose is to help in proving a theorem. It is a stepping stone on the path to proving a theorem. Very occasionally lemmas can take on a life of their own (Zorn's lemma, Urysohn's lemma, Burnside's lemma, Sperner's lemma).
- (4) **Corollary**—a result in which the (usually short) proof relies heavily on a given theorem (we often say that "this is a corollary of Theorem A").
- (5) **Proposition**—a proved and often interesting result, but generally less important than a theorem.
- (6) **Conjecture**—a statement that is unproved, but is believed to be true (Collatz conjecture, Goldbach conjecture, twin prime conjecture).
- (7) **Claim**—an assertion that is then proved. It is often used like an informal lemma.
- (8) **Axiom/Postulate**—a statement that is assumed to be true without proof. These are the basic building blocks from which all theorems are proved (Euclid's five postulates, Zermelo-Frankel axioms, Peano axioms).
- (9) **Identity**—a mathematical expression giving the equality of two (often variable)

Approximate Policy Evaluation

- In practice, we can't perform an infinite number of iterations.
- Suppose that we perform value iteration for n steps and $\|H^{(n)}(V) - H^{(n-1)}(V)\|_{\infty} = \epsilon$, how far is $H^{(n)}(V)$ from V^{π} ?

Approximate Policy Evaluation

- **Theorem 3:** If $\|H^{(n)}(V) - H^{(n-1)}(V)\|_{\infty} \leq \epsilon$ then

$$\|V^{\pi} - H^{(n)}(V)\|_{\infty} \leq \frac{\epsilon}{1 - \gamma}$$

- **Proof** $\|V^{\pi} - H^{(n)}(V)\|_{\infty}$

$$= \|H^{(\infty)}(V) - H^{(n)}(V)\|_{\infty} \quad (\text{by Theorem 2})$$

$$= \left\| \sum_{t=1}^{\infty} H^{(t+n)}(V) - H^{(t+n-1)}(V) \right\|_{\infty}$$

$$\leq \sum_{t=1}^{\infty} \|H^{(t+n)}(V) - H^{(t+n-1)}(V)\|_{\infty} \quad (\|A + B\| \leq \|A\| + \|B\|)$$

$$= \sum_{t=1}^{\infty} \gamma^t \epsilon = \frac{\epsilon}{1 - \gamma} \quad (\text{by Lemma 1})$$

Optimal Value Function

- Non-linear system of equations

$$V_{\infty}^*(s) = \max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V_{\infty}^*(s') \quad \forall s$$

- Matrix form:

R^a : $|S| \times 1$ column vector of rewards for a

V^* : $|S| \times 1$ column vector of optimal values

T^a : $|S| \times |S|$ matrix of transition prob for a

$$V^* = \max_a R^a + \gamma T^a V^*$$

Contraction

- Let $H^*(V) \stackrel{\text{def}}{=} \max_a R^a + \gamma T^a V$ be the operator in value iteration
- **Lemma 4:** H^* is a **contraction mapping**.

$$\|H^*(\tilde{V}) - H^*(V)\|_\infty \leq \gamma \|\tilde{V} - V\|_\infty$$

- Proof: without loss of generality,

let $H^*(\tilde{V})(s) \geq H^*(V)(s)$ and

let $a_s^* = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s')$

$\tilde{a}_s^* = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) \tilde{V}(s')$

Contraction

- Proof continued:
- Then $0 \leq H(\tilde{V})(s) - H^*(V)(s)$ (by assumption)
 $= R(s, \tilde{a}_s^*) + \gamma \sum_{s'} \Pr(s'|s, \tilde{a}_s^*) \tilde{V}(s')$ (by definition)
 $\quad - R(s, a_s^*) + \gamma \sum_{s'} \Pr(s'|s, a_s^*) V(s')$
 $\leq R(s, \tilde{a}_s^*) + \gamma \sum_{s'} \Pr(s'|s, \tilde{a}_s^*) \tilde{V}(s')$ (since \tilde{a}_s^* suboptimal for V)
 $\quad - R(s, \tilde{a}_s^*) + \gamma \sum_{s'} \Pr(s'|s, \tilde{a}_s^*) V(s')$
 $= \gamma \sum_{s'} \Pr(s'|s, \tilde{a}_s^*) [\tilde{V}(s') - V(s')]$
 $\leq \gamma \sum_{s'} \Pr(s'|s, \tilde{a}_s^*) \|\tilde{V} - V\|_\infty$ (maxnorm upper bound)
 $= \gamma \|\tilde{V} - V\|_\infty$ (since $\sum_{s'} \Pr(s'|s, \tilde{a}_s^*) = 1$)
- Repeat the same argument for $H^*(V)(s) \geq H^*(\tilde{V})(s)$ and for each s

Convergence

- **Theorem 5:** Value iteration converges to V^* for any initial estimate V

$$\lim_{n \rightarrow \infty} H^{*(n)}(V) = V^* \quad \forall V$$

- Proof
 - By definition $V^* = H^{*(\infty)}(0)$, but value iteration computes $H^{*(\infty)}(V)$ for some initial V
 - By Lemma 4, $\|H^{*(n)}(V) - H^{*(n)}(\tilde{V})\|_{\infty} \leq \gamma^n \|V - \tilde{V}\|_{\infty}$
 - Hence, when $n \rightarrow \infty$, then $\|H^{*(n)}(V) - H^{*(n)}(0)\|_{\infty} \rightarrow 0$ and $H^{*(\infty)}(V) = V^* \quad \forall V$

Value Iteration

- Even when horizon is infinite, perform finitely many iterations
- Stop when $\|V_n - V_{n-1}\| \leq \epsilon$

valueiteration(MDP)

$V_0^* \leftarrow \max_a R^a; \quad n \leftarrow 0$

Repeat

$n \leftarrow n + 1$

$V_n \leftarrow \max_a R^a + \gamma T^a V_{n-1}$

Until $\|V_n - V_{n-1}\|_\infty \leq \epsilon$

Return V_n

Induced Policy

- Since $\|V_n - V_{n-1}\|_\infty \leq \epsilon$, by Theorem 5: we know that $\|V_n - V^*\|_\infty \leq \frac{\epsilon}{1-\gamma}$
- But, how good is the stationary policy $\pi_n(s)$ extracted based on V_n ?

$$\pi_n(s) = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V_n(s')$$

- How far is V^{π_n} from V^* ?

Induced Policy

- **Theorem 6:** $\|V^{\pi_n} - V^*\|_{\infty} \leq \frac{2\epsilon}{1-\gamma}$

- **Proof**

$$\begin{aligned}\|V^{\pi_n} - V^*\|_{\infty} &= \|V^{\pi_n} - V_n + V_n - V^*\|_{\infty} \\ &\leq \|V^{\pi_n} - V_n\|_{\infty} + \|V_n - V^*\|_{\infty} \quad (\|A + B\| \leq \|A\| + \|B\|) \\ &= \|H^{\pi_n(\infty)}(V_n) - V_n\|_{\infty} + \|V_n - H^{*(\infty)}(V_n)\|_{\infty} \\ &\leq \frac{\epsilon}{1-\gamma} + \frac{\epsilon}{1-\gamma} \quad (\text{by Theorems 3 and 5}) \\ &= \frac{2\epsilon}{1-\gamma}\end{aligned}$$

Summary

- Value iteration
 - Simple dynamic programming algorithm
 - Complexity: $O(n|A||S|^2)$ (see P16)
 - Here n is the number of iterations
- Can we optimize the policy directly instead of optimizing the value function and then inducing a policy?
 - Yes: by policy iteration