

4

Last Time

Value Function - Expectation of future rewards starting at s

$$V^*(s) = \max_a \left(R(s,a) + \gamma E[V^*(s')] \right)$$

Two Algorithms

Policy Iteration

Policy Evaluation

V^π

Policy Improvement

Value Iteration

$$V' = B[V]$$

$$B[V](s) = \max_a \left(R(s,a) + \gamma E[V(s')] \right)$$

$$\sum_{s'} T(s'|s,a) V(s')$$

Today

V.I. converges

V^* unique

Julia perf + debugging

Q value

$$Q^\pi(s,a) = R(s,a) + \gamma E[V^\pi(s')]$$

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

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s,a) \leftarrow$$

~~$\operatorname{argmax}_a R(s,a)$~~ Myopic

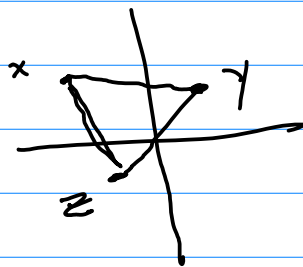
$$V_{k+1} = B[V_k]$$

$M = \mathbb{R}^{151}$

Theorem Let $\{V_k\}_{k=1}^{\infty}$ be a sequence of value fns for a discrete MDP, generated by applying B . If $\gamma < 1$, then $\lim_{k \rightarrow \infty} V_k = V^*$.

Def Let M be a set. A metric on M is a function $d: M \times M \rightarrow [0, \infty)$ that satisfies

- i) $d(x, y) = 0$ iff $x = y$
- ii) $d(x, y) = d(y, x) \quad \forall x, y \in M$
- iii) $d(x, z) \leq d(x, y) + d(y, z) \quad \forall x, y, z \in M$



Def A contraction mapping on set M is a function $f: M \rightarrow M$ satisfying

$$d(f(x), f(y)) \leq c d(x, y)$$

for some c , $0 \leq c < 1$ and all $x, y \in M$

Def x^* is said to be a fixed point of f if $f(x^*) = x^*$

Banach's Theorem If f is a contraction mapping on (M, d) , then

- i) f has a single, unique fixed point
- ii) If $\{x_n\}$ is a sequence defined by $x_{k+1} = f(x_k)$ then $\lim_{k \rightarrow \infty} x_k = x^*$

$$\begin{aligned}\|V\|_1 &= \sum V \\ \|V\|_2 &= \sqrt{\sum V_i^2} \\ \|V\|_\infty &= \max_i |V_i|\end{aligned}$$

Prove that

$$\|V_1 - V_2\|_\infty \text{ is a metric}$$

$$(i) \|V_1 - V_2\|_\infty = \|V_2 - V_1\|_\infty$$

$$\begin{aligned}(ii) \quad & \max_i |x_i - z_i| = \max_i |x - y + y - z| \quad \leftarrow \\ & \left[\max_i |x_i - y_i + y_i - z_i| \leq |x_i - y_i| + |y_i - z_i| \right] \quad \leftarrow \text{triangle inequality} \\ & \max |x - y + y - z| \leq \max (|x - y| + |y - z|) \\ & \leq \max |x - y| + \max |y - z| \\ & \quad \underline{\quad \quad \quad} d(x, y) + d(y, z)\end{aligned}$$

Lemma 1 $\|V_1 - V_2\|_\infty$ is a metric on $\mathbb{R}^{|S|}$

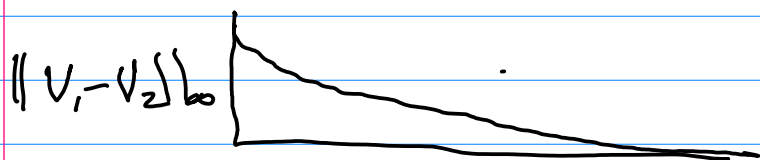
$$\begin{aligned}\|B[V_1] - B[V_2]\|_\infty &= \max_{s \in S} |B[V_1](s) - B[V_2](s)| \\ &= \max_{s \in S} \left| \max_{a \in A} (R(s, a) + \gamma \sum_{s'} T(s'|s, a) V_1(s')) - \max_{a \in A} (R(s, a) + \gamma \sum_{s'} T(s'|s, a) V_2(s')) \right| \\ &\leq \max_{s \in S} \left| \max_{a \in A} \left(R(s, a) + \gamma \sum_{s'} T(s'|s, a) V_1(s') \right) - \left(R(s, a) + \gamma \sum_{s'} T(s'|s, a) V_2(s') \right) \right| \\ &= \max_{s \in S} \max_{a \in A} \left| \sum_{s'} T(s'|s, a) (V_1(s') - V_2(s')) \right| \\ &\leq \max_{s \in S} \max_{a \in A} \sum_{s'} T(s'|s, a) |V_1(s') - V_2(s')| \\ &\leq \max_{s \in S} \max_{a \in A} \sum_{s'} T(s'|s, a) \|V_1 - V_2\|_\infty\end{aligned}$$

$$= \gamma \|V_1 - V_2\|_\infty \max_a \sum_{s'} T(s'|s, a)$$

$$\leq \gamma \|V_1 - V_2\|_\infty$$

Lemma 2 B is a contraction mapping on $\mathbb{R}^{|S|}$

By L1, L2, Banach's Theorem, we have proved the original Theorem \square



$S = \{1, 2, 3\}$ $|S| = 3$ "cardinality" size of S

$\mathbb{R}^{|S|}$ = 3-dimensional vector space

$$V_1 = [0.0, 0.1, 0.2]$$

$$V_2 = [0.1, 0.3, 0.5]$$

Breakout Rooms

$$S = \{1, 2, 3\}$$

$$A = \{1, -1\}$$

$$V^*(3) = 98$$

$$R(2, a) = 0$$

$$R(1, a) = 1$$

$$\gamma = 0.99$$

$$T(s'|s, a) = \begin{cases} 0.7 & \text{if } s' = \text{clamp}(s+a, 1, 3) \\ 0.3 & \text{if } s' = \text{clamp}(s-a, 1, 3) \\ 0 & \text{o.w. } R_{\text{opposite}} \end{cases}$$

desired

Optimal policy? $\pi^*(s) \rightarrow$

$$V_1 = 95$$

$$V_2 = 95$$

$$R(s, a) + \gamma E(V(s'))$$

$$V_1 = 1 + \gamma(0.7V_1 + 0.3V_2)$$

$$V_2 = 0 + \gamma(0.7V_1 + 0.3 \cdot 98)$$

$$V_1 = 96.67$$

$$V_2 = 96.62$$

$$1 + \gamma(0.7V_2 + 0.3V_1)$$

$$0 + \gamma(0.7 \cdot 98 + 0.3V_1)$$