# Module 6
# Value Iteration

CS 886 Sequential Decision Making and
Reinforcement Learning

University of Waterloo

# Markov Decision Process

- Definition
  - Set of states: $S$
  - Set of actions (i.e., decisions): $A$
  - Transition model: $\Pr(s_t | s_{t-1}, a_{t-1})$
  - Reward model (i.e., utility): $R(s_t, a_t)$
  - Discount factor: $0 \leq \gamma \leq 1$
  - Horizon (i.e., # of time steps): $h$

- Goal: find optimal policy $\pi$

2

# Finite Horizon

- Policy evaluation

$$V_h^\pi(s) = \sum_{t=0}^{h} \gamma^t \Pr(S_t = s' | S_0 = s, \pi) R(s', \pi_t(s'))$$

- Recursive form (dynamic programming)

$$V_0^\pi(s) = R(s, \pi_0(s))$$

$$V_t^\pi(s) = R\big(s, \pi_t(s)\big) + \gamma \sum_{s'} \Pr(s' | s, \pi_t(s)) \, V_{t-1}^\pi(s')$$

# Finite Horizon

- Optimal Policy $\pi^*$

$$V_h^{\pi^*}(s) \geq V_h^{\pi}(s) \quad \forall \pi, s$$

- Optimal value function $V^*$ (shorthand for $V^{\pi^*}$)

$$V_0^*(s) = \max_a R(s, a)$$

$$\underbrace{V_t^*(s) = \max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V_{t-1}^*(s')}_{\text{Bellman's equation}}$$

4

# Value Iteration Algorithm

valueIteration(MDP)
$V_0^*(s) \leftarrow \max_a R(s,a) \ \forall s$
For $t = 1$ to $h$ do
$\qquad V_t^*(s) \leftarrow \max_a R(s,a) + \gamma \sum_{s'} \Pr(s'|s,a) V_{t-1}^*(s') \ \forall s$
Return $V^*$

Optimal policy $\pi^*$

$t = 0:\ \pi_0^*(s) \leftarrow \underset{a}{\text{argmax}}\, R(s,a) \ \forall s$

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

NB: $\pi^*$ 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

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

<span style="color:darkred">valueIteration(MDP)</span>

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

For $t = 1$ to $h$ do

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

Return $V^*$

6

# Infinite Horizon

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

- Policy evaluation:

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

- Bellman's equation:

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

# Policy evaluation

- Linear system of equations

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

- Matrix form:

  $R$: $|S| \times 1$ column vector of sate rewards for $\pi$

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

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

$$V = R + \gamma T V$$

# 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 $V \leftarrow R + \gamma TV$

9

# Contraction

- Let $H(V) \stackrel{\text{def}}{=} R + \gamma TV$ be the policy eval operator

- Lemma 1: $H$ is a contraction mapping.

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

- Proof $\left\| H(\tilde{V}) - H(V) \right\|_{\infty}$

$$= \left\| R + \gamma T\tilde{V} - R - \gamma TV \right\|_{\infty} \quad \text{(by definition)}$$

$$= \left\| \gamma T(\tilde{V} - V) \right\|_{\infty} \quad \text{(simplification)}$$

$$\leq \gamma \left\| T \right\|_{\infty} \left\| \tilde{V} - V \right\|_{\infty} \quad \text{(since } \|AB\| \leq \|A\|\|B\| \text{)}$$

$$= \gamma \left\| \tilde{V} - V \right\|_{\infty} \quad \text{(since } \max_{s} \sum_{s'} T(s, s') = 1 \text{)}$$

10

# Convergence

- Theorem 2: <span style="color:darkred">Policy evaluation converges to $V^\pi$</span> for any initial estimate $V$

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

- Proof

  - By definition $V^\pi = H^{(\infty)}(0)$, but policy evaluation computes $H^{(\infty)}(V)$ for any initial $V$

  - By lemma 1, $\left|\left|H^{(n)}(V) - H^{(n)}(\tilde{V})\right|\right|_\infty \leq \gamma^n \left|\left|V - \tilde{V}\right|\right|_\infty$

  - Hence, when $n \to \infty$, then $\left|\left|H^{(n)}(V) - H^{(n)}(0)\right|\right|_\infty \to 0$ and $H^{(\infty)}(V) = V^\pi \quad \forall V$

# Approximate Policy Evaluation

- In practice, we can't perform an infinite number of iterations.

- Suppose that we perform value iteration for $k$ steps and $\left\| H^{(k)}(V) - H^{(k-1)}(V) \right\|_\infty = \epsilon$, how far is $H^{(k)}(V)$ from $V^\pi$?

# Approximate Policy Evaluation

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

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

- Proof $\left\|V^\pi - H^{(k)}(V)\right\|_\infty$

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

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

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

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

13

# Optimal Value Function

- Non-linear system of equations

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

- 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 $a$

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

14

# Contraction

- Let $H^*(V) \overset{\text{def}}{=} \max_a R^a + \gamma T^a V$ be the operator in value iteration

- Lemma 3: $H^*$ is a contraction mapping.

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

- Proof: without loss of generality,
  let $H^*(\tilde{V})(s) \geq H^*(V)(s)$ and
  let $a_s^* = \underset{a}{\text{argmax}}\, R(s,a) + \gamma \sum_{s'} \Pr(s'|s,a)\, V(s')$

15

# Contraction

- Proof continued:

- Then $0 \leq H^*(\tilde{V})(s) - H^*(V)(s)$       (by assumption)

$$\leq R(s, a_s^*) + \gamma \sum_{s'} \Pr(s'|s, a_s^*) \tilde{V}(s') \quad \text{(by definition)}$$
$$-R(s, a_s^*) - \gamma \sum_{s'} \Pr(s'|s, a_s^*) V(s')$$
$$= \gamma \sum_{s'} \Pr(s'|s, a_s^*) [\tilde{V}(s') - V(s')]$$
$$\leq \gamma \sum_{s'} \Pr(s'|s, a_s^*) \left|\left|\tilde{V} - V\right|\right|_\infty \quad \text{(maxnorm upper bound)}$$
$$= \gamma \left|\left|\tilde{V} - V\right|\right|_\infty \quad \text{(since } \sum_{s'} \Pr(s'|s, a_s^*) = 1\text{)}$$

- Repeat the same argument for $H^*(V)(s) \geq H^*(\tilde{V})(s)$ and for each $s$

16

# Convergence

- Theorem 4: Value iteration converges to $V^*$ for any initial estimate $V$

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

- Proof

  - By definition $V^* = H^{*(\infty)}(0)$, but value iteration computes $H^{*(\infty)}(V)$ for some initial $V$

  - By lemma 3, $\left\| H^{*(n)}(V) - H^{*(n)}(\tilde{V}) \right\|_\infty \leq \gamma^n \left\| V - \tilde{V} \right\|_\infty$

  - Hence, when $n \to \infty$, then $\left\| H^{*(n)}(V) - H^{*(n)}(0) \right\|_\infty \to 0$ 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)

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

Repeat

$\qquad n \leftarrow n + 1$

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

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

Return $V_n$

18

# Induced Policy

- Since $\left\|V_n - V_{n-1}\right\|_\infty \leq \epsilon$, by Theorem 4: we know that $\left\|V_n - V^*\right\|_\infty \leq \frac{\epsilon}{1-\gamma}$

- 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 $V^{\pi_n}$ from $V^*$?

# Induced Policy

- Theorem 5: $\left\|V^{\pi_n} - V^*\right\|_\infty \leq \frac{2\epsilon}{1-\gamma}$

- Proof

$$\left\|V^{\pi_n} - V^*\right\|_\infty = \left\|V^{\pi_n} - V_n + V_n - V^*\right\|_\infty$$

$$\leq \left\|V^{\pi_n} - V_n\right\|_\infty + \left\|V_n - V^*\right\|_\infty \quad (\|A+B\| \leq \|A\| + \|B\|)$$

$$= \left\|H^{\pi_n(\infty)}(V_n) - V_n\right\|_\infty + \left\|V_n - H^{*(\infty)}(V_n)\right\|_\infty$$

$$\leq \frac{\epsilon}{1-\gamma} + \frac{\epsilon}{1-\gamma} \quad \text{(by Theorems 2 and 4)}$$

$$= \frac{2\epsilon}{1-\gamma}$$

20

# Summary

- Value iteration
    - Simple dynamic programming algorithm
    - Complexity: $O(n|A||S|^2)$
        - 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

21