

## Last Time

MDP ( $s, A, T, R, \gamma, [p_0]$ )

Policy:  $\pi: S \rightarrow A$        $a = \pi(s)$

MC Policy Eval.

---

MDP: Balance immediate and future rewards

---

## Today

- When is enough MC sims? (SEM)
- Policy Search (Cross Entropy)
- Value Function
- Policy Iter.      } Optimal Policies,
- Value Iter.      } "Dynamic Programming"

## MC Evaluation

$$\hat{u} = \sum_{t=0}^{T-1} \gamma^t r_t$$

$$\bar{u}_m = \frac{1}{m} \sum_{k=1}^m \hat{u}_k$$

$$\begin{aligned} \sum_{t=0}^{\infty} \gamma^t r_t &= \sum_{t=0}^{T-1} \gamma^t r_t + \underbrace{\sum_{t=T}^{\infty} \gamma^t r_t}_{\leq \gamma^T \frac{r}{1-\gamma}} \\ &= \sum_{t=0}^{T-1} \gamma^t r_t + \gamma^T \sum_{t=T}^{\infty} \gamma^{t-T} r_t \end{aligned}$$

Central Limit Theorem

as  $m \rightarrow \infty$   $\bar{u}_m \xrightarrow{d} N\left(\overset{\text{true mean return}}{\mu}, \overset{\text{std. dev of } \hat{u}}{\frac{\sigma_{\hat{u}}}{\sqrt{m}}}\right) \leq \gamma^T \frac{r}{1-\gamma}$

$$\boxed{\text{SEM}(\bar{u}_m) \equiv \frac{\text{std}(\hat{u})}{\sqrt{m}}}$$

## Policy Search (Ch 10)

- ① Parameterize policy  $\pi(s) = f(s; \theta)$
- ② Optimize  $\theta$  w/ favorite 0th order optimization alg.  $\pi(s) = -k s$   
 $\theta$   
 $f_{\min \text{ con}}$

$$|S| < \infty \quad |A| < \infty$$

Value Function

$$\pi(s) = \underset{a \in A}{\operatorname{argmax}} R(s, a)$$

Not good

"Myopic"

"Greedy"

$$V^{\pi}(s) = E\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_t = \pi(s_t)\right]$$

Hand

$$V^\pi(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_t = \pi(s_t) \right]$$

$$= R(s, \pi(s)) + E \left[ \sum_{t=1}^{\infty} \gamma^t r_t \mid s_1 \sim T(s_1), a_t = \pi(s_t) \right]$$

$$= \quad \quad + \gamma E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s', a_t = \pi(s_t) \right]$$

$$\rightarrow V^\pi(s) = R(s, \pi(s)) + \gamma E_{s' \sim T(s, \pi)} [V^\pi(s')] ]$$

$$= \quad \quad + \gamma \sum_{s' \in S} T(s'|s, a) V^\pi(s') \quad \left\{ \begin{array}{l} T_{i,j}^\pi = T(j|i, \pi(i)) \\ R_i^\pi = R(i, \pi(i)) \end{array} \right.$$

$$\rightarrow V^\pi = R^\pi + \gamma T^\pi V^\pi$$

$$V^\pi - \gamma T^\pi V^\pi = R^\pi$$

$$V^\pi = (I - \gamma T^\pi)^{-1} R^\pi$$

matrix policy evaluation

## Policy Iteration

initialize  $\pi, \pi'$

while  $\pi \neq \pi'$

$\pi = \pi'$

$$\rightarrow V^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi \quad \swarrow \text{policy eval}$$

$$\pi'(s) \leftarrow \underset{a \in A}{\operatorname{argmax}} \left( R(s, a) + \gamma \sum_{s'} T(s'|s, a) V^\pi(s') \right)$$

$\forall s \in S$

$\pi$  is optimal!

$\uparrow$  policy update

know  $V^*$

①

$$\pi^*(s) = \underset{a \in A}{\operatorname{argmax}} R(s, a) + \gamma E_{s' \sim T(s, a)} [V^*(s')] ]$$

if know  $V^*$ , we know  $\pi^*$

## Bellman's Equation

②

$$V^*(s) = \max_{a \in A} \left[ R(s, a) + \gamma \mathbb{E}_{s' \sim T(s, a)} [V^*(s')] \right]$$

### Value Iteration

Why does  
VI converge?

initialize  $V, V'$

while  $V \neq V'$

$V = V'$

$$V'(s) \leftarrow \max_a R(s, a) + \gamma \mathbb{E}_{s' \sim T(s, a)} [V(s')]$$

$\forall s \in S$

when terminates,  $V = V^*$ .

Fact: There can be many optimal policies,  
but only one optimal value function!

"Dynamic Programming"

↑ Because it sounds cool

Break down a large opt. into many smaller ones.