# Solving MDPs

# Value Iteration

Repeated Bellman updates:

---

1: **for all** States $s$ **do**

2:    $U(s) \leftarrow R(s)$

3: **while** unsatisfied **do**

4:    **for each** state $s$ **do**

5:       $U'(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$

6:    $U \leftarrow U'$

---

Utility values are guaranteed to converge with enough updates
This equilibrium gives an optimal policy

# When To Stop

Solving MDPs
Value Iteration
Stopping
Policy Iteration

Reinforcement
Learning

Eventually updates have diminishing returns

$$||U_{i+1} - U_i|| = \text{max difference of } U \text{ between iterations}$$

$$||U_{i+1} - U_i|| < \varepsilon(1 - \gamma)/\gamma \implies ||U_{i+1} - U^*|| < \varepsilon$$

**Policy loss**: the most utility an agent loses by following policy $\pi_i$ instead of $\pi^*$

Policy is often optimal many iterations before $U_i$ converges on optimal

Solving MDPs
Value Iteration
Stopping
Policy Iteration

Reinforcement
Learning

# Policy Iteration

$U$: utility for all states

$\pi$: action policy for all states

---

POLICYITER($mdp$)

1: **repeat**
2:   $U \leftarrow$ PolicyEval($\pi, U, mdp$)
3:   $unchanged \leftarrow$ true
4:   **for all** States $s$ **do**
5:     $a^* \leftarrow \operatorname{argmax}_{a \in A(s)} QVal(mdp, s, a, U)$
6:     **if** $QVal(mdp, s, a^*, U) > QVal(mdp, s, \pi[s], U)$ **then**
7:       $\pi[s] \leftarrow a^*$
8:       $unchanged \leftarrow$ false
9: **until** $unchanged$
10: **return** $\pi$

---

# Reinforcement Learning

Solving MDPs

Reinforcement
Learning

Definition

ADP

Sweeping

# Definition

Build a good policy based on experience only: $(s, a, s', r)$
Objective:

- **finite horizon**: $R(s_0) + R(s_1) + R(s_2)$, or
- **infinite horizon**: $R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \ldots$

*Learning* because we don't start with a model of the system

# Adaptive Dynamic Programming

- Learn $T$ and $R$ functions through trials
- Easy approach to find rewards from each state, and transition probabilities
- Calculate $\pi$ using MDP solution

Update utility and policy of all states until satisfied

# Prioritized Sweeping

Focus on areas of model where large changes are expected

---

$\text{UPDATE}(s, a, s', r)$

---

1: update model
2: $\text{UPDATE}(s)$
3: Do $k$ times: update highest priority

---

$\text{UPDATE}(s)$

---

1: update $U(s)$
2: priority of s $\leftarrow 0$
3:
4: **for all** States $s'$ that are predecessors of $s$ **do**
5:     priority of $s' \leftarrow \max(current, \max_a \delta \hat{T}(s', a, s))$

---