# What is an MDP?

Optimization Problem         maximize $E\left[\sum_{t=0}^{\infty} R(s_t, a_t)\right]$

Defined by     $(S, A, T, R, \gamma)$

# What is a policy?

- Closed loop, <u>deterministic</u>         $\pi: S \to A$
  "policy"

- Open loop, <u>deterministic</u>:     List of actions executed
                                                          in order

For deterministic problems         Finite  Horizon $T$
$\exists$ optimal open loop policy      O.L.                    C.L.
                                                      $|A|^T$               $|A|^{|S| T}$

## Today

What is a Value Function?
How can we find optimal policies?
    Dynamic Programming
    Policy Iteration
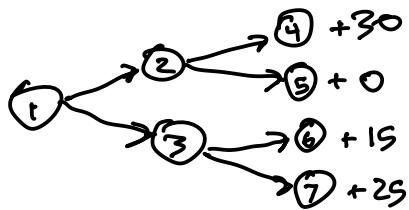
---

$\pi(s) = \arg\max_{a \in A} R(s, a)$         "Myopic" / "Greedy"
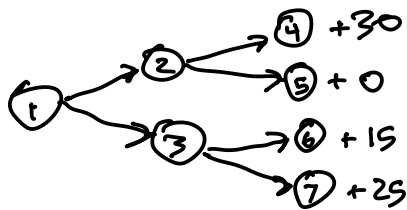


+1   for $\to$
−10   for 🔥

# Dynamic Programming

Large problem → smaller sub

U/D

## Naive

Consider all possible plans

$$R(U) + R(U) = 30 \leftarrow$$
$$U \qquad D = 0$$
$$D \qquad U = 15$$
$$D \qquad D = 25$$

## Value / Utility Function

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



### DP

$$V(4) = 30$$
$$V(5) = 0$$
$$V(6) = 15$$
$$V(7) = 25$$

3 computation steps
$$\begin{cases} \text{at } 2 & \pi^*(2) = U \qquad V^*(2) = 30 \\ \text{at } 3 & \pi^*(3) = D \qquad V^*(3) = 25 \\ \text{at } 1 & \\ & \qquad \pi^*(1) = U \quad V^*(1) = 30 \end{cases}$$

## Two Basic Algorithms

Policy Iteration ← Easier to Understand

Value Iteration ← Easier to Implement

Bellman's principle of optimality

Every sub-path of an optimal path is optimal

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

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

$$= R(s, \pi(s)) + \gamma E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 \sim T(s,a), a_t = \pi(s_t)\right]$$

$$V^{\pi}(s) = R(s, \pi(s)) + \gamma E\left[V^{\pi}(s') \mid s' \sim T(s,a)\right]$$

Discrete state/action spaces

$$E_{s' \sim T(s,p)}\left[V^{\pi}(s')\right]$$

$$V^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} T(s' \mid s, \pi(s)) V^{\pi}(s')$$

$$V^{\pi} = R^{\pi} + \gamma T^{\pi} V^{\pi} \quad \leftarrow |S| \text{ vector}$$
$|S|$ vector $\quad |S|$ vector $\quad |S| \times |S|$ matrix

$$T^{\pi}_{ij} = T^{\pi}(s' = j \mid s = i, \pi(i))$$

$|S| = $ number of element in $S$

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

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

$$V^{\pi} = (I - \gamma T^{\pi})^{-1} R^{\pi} \quad \leftarrow \text{exact policy evaluation}$$

$$O(|S|^3) \quad \text{time}$$

Policy iteration
$\pi_0 = $ guess, $k = 0$
while $\pi_k \neq \pi_{k+1}$

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

$$Q^{\pi}(s,a)$$

for $s \in S$

$$\pi_{k+1}(s) = \underset{a \in A}{\operatorname{argmax}} \left(R(s,a) + \gamma \sum_{s' \in S} T(s' \mid s,a) V^{\pi_k}(s')\right)$$

$k = k+1$

return $\pi_k$

optim

| +10 | | +10 |
|---|---|---|
| 1 | 2 | 3 |

$\pi_a(2) = \leftarrow \quad V^{\pi_a}(2) = \gamma 10$
$\pi_b(2) = \rightarrow \quad V^{\pi_b}(2) = \gamma 10$

| +8 | | +9 |
|---|---|---|

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

Bellman's Equation

$$V^{k+1}(s) = \max_{a \in A} \left(R(s,a) + \gamma E\left[V^k(s') \mid s' \sim T(s,a)\right]\right) \quad \text{repeat}$$