

Last Time

$$(S, A, \underset{\uparrow}{T}, R, \gamma)$$

Conditional
Probability Dist.

Last Time

- How is a **Markov decision process** defined?

Last Time

- How is a **Markov decision process** defined?
- What is a **policy**?

$$\pi: S \rightarrow A$$

Last Time

- How is a **Markov decision process** defined?
- What is a **policy**?
- How do we **evaluate** policies?

$$E. \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$
$$\sum_{k=0}^S \frac{1}{S} \sum_{t=0}^T \gamma^t r_t$$

Guiding Questions

Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

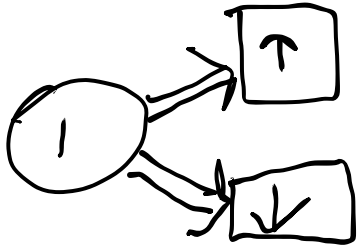
MDP Example: Up-Down Problem

1

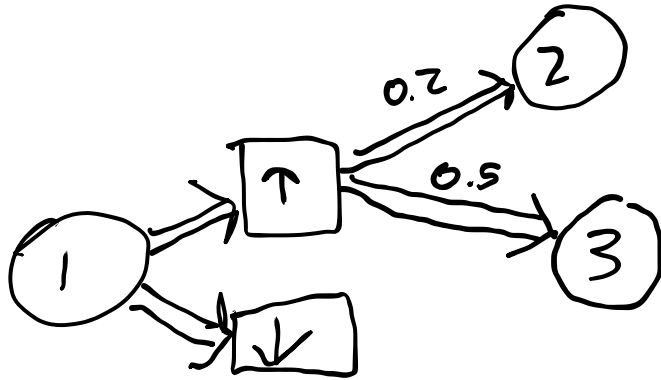
MDP Example: Up-Down Problem



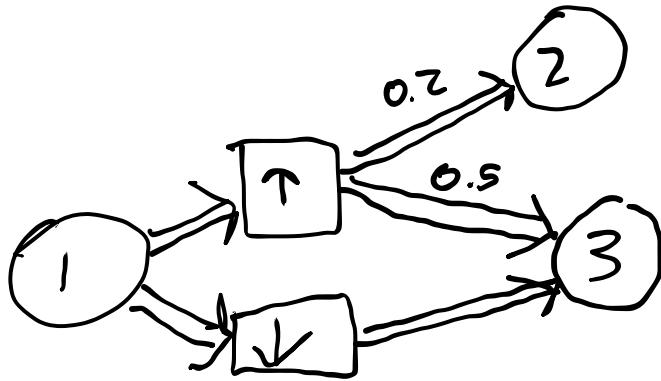
MDP Example: Up-Down Problem



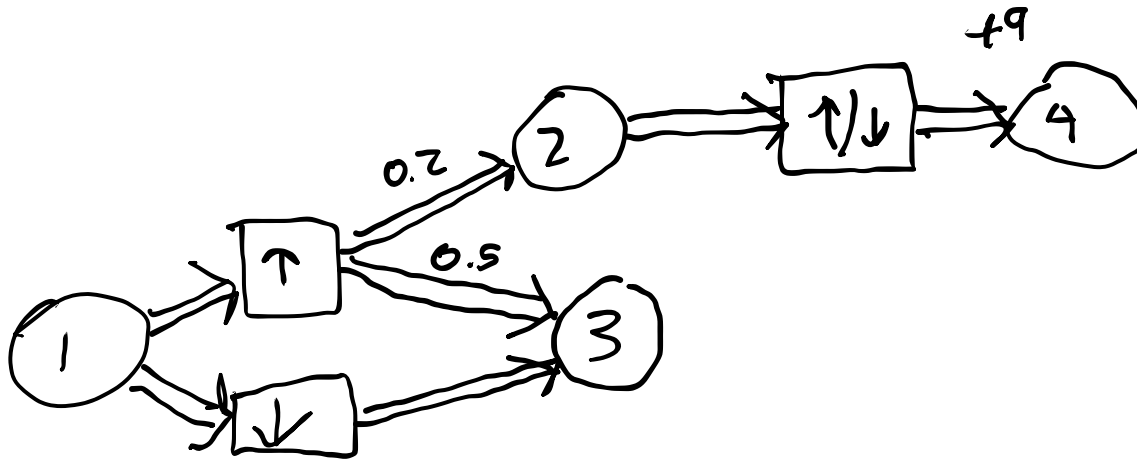
MDP Example: Up-Down Problem



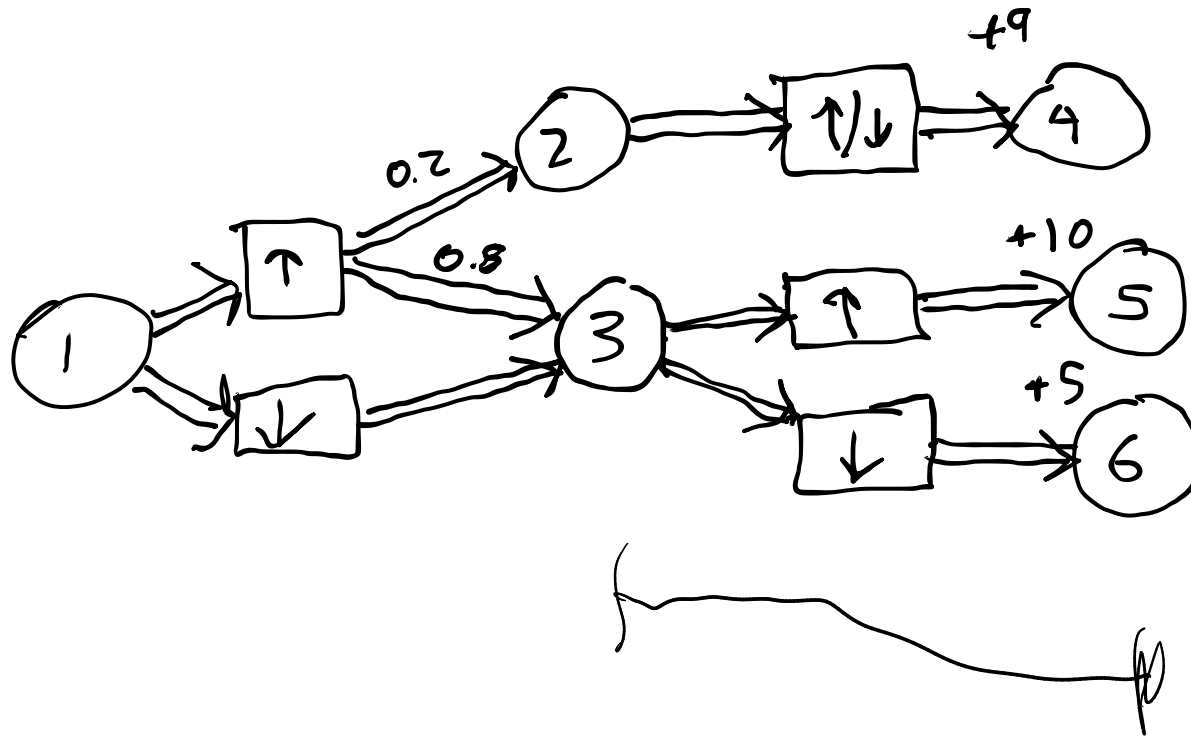
MDP Example: Up-Down Problem



MDP Example: Up-Down Problem



MDP Example: Up-Down Problem



s	$\pi(s)$
1	↓
3	↑

Value Functions

$$E[r_0 | s_0 = s] = R(s, \pi(s))$$

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

myopic $\pi(s) = \operatorname{argmax}_a R(s, a)$

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

$$= \gamma^0 R(s, \pi(s)) + E \left[\sum_{t=1}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

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

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

$$+ \gamma E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s' \right]$$

$s' \sim T(s, a)$

$$V^\pi(s) = \underbrace{R(s, \pi(s))}_{\text{Immediate reward}} + \gamma \underbrace{E[V^\pi(s')]_{s' \sim T(s, a)}}_{\text{Rewards in future}}$$

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

$V = U$
in book

Exact ^{policy} _V Matrix Evaluation

$$V^\pi = [v \dots v]$$

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

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

$$V_i^\pi = R_i^\pi + \gamma \sum_{j \in \{1, \dots, |S|\}} T_{ij}^\pi V_j^\pi$$

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

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

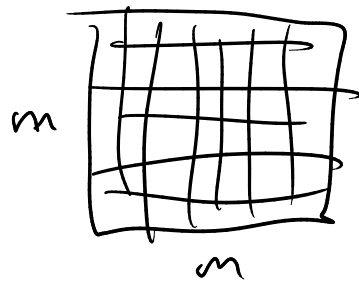
$$(I - \gamma T^\pi) V^\pi = R^\pi$$

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

iterative
policy
evaluation

$$T_{ij} = T(s'=j | s=i, a=\pi(i))$$

$$O(n^3)$$



$$|S| \times |S|$$

$$O(s^6)$$

$$|S| = m^d$$

Policy Iteration

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

initialize π, π' (differently)

while $\pi \neq \pi'$

$\pi \leftarrow \pi'$

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

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

return $\pi \leftarrow \text{optimal}$

policy evaluation

policy update

$Q^\pi(s, a)$

Bellman's Equation

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

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

$$V^*$$

$$V^*(s) = R(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim T(s, \pi^*(s))} [V^*(s')]$$

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

Value Function Policies

$$\pi^*(s) = \operatorname{argmax}_a \left\{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)} [V^*(s')] \right\}$$

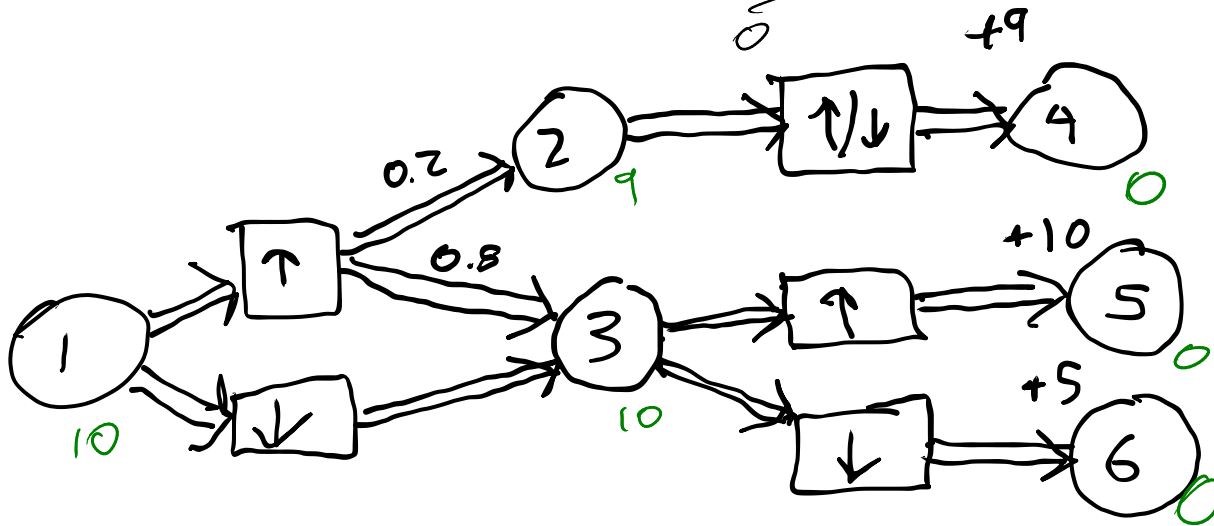
$$V^* \Rightarrow \pi^*$$

$$\pi^* \Rightarrow V^*$$

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

Backup by hand example

$$V^*(s) = \max_a \{ R(s,a) + \gamma \underbrace{E[V^*(s')]}_0 \}$$



s	$V^*(s)$	a	$R(s,a)$	$E[V^*(s')]$
6	0			
5	0			
4	0			
3	10	↑	+10	0
3		↓	+5	0
2	9	↑	+9	0
2		↓	+9	0
1	10	↑	0	$0.2 \cdot 9 + 0.8 \cdot 10 = 9.8$
1		↓	0	10

Breakout Rooms: DIA Run

Boulder.

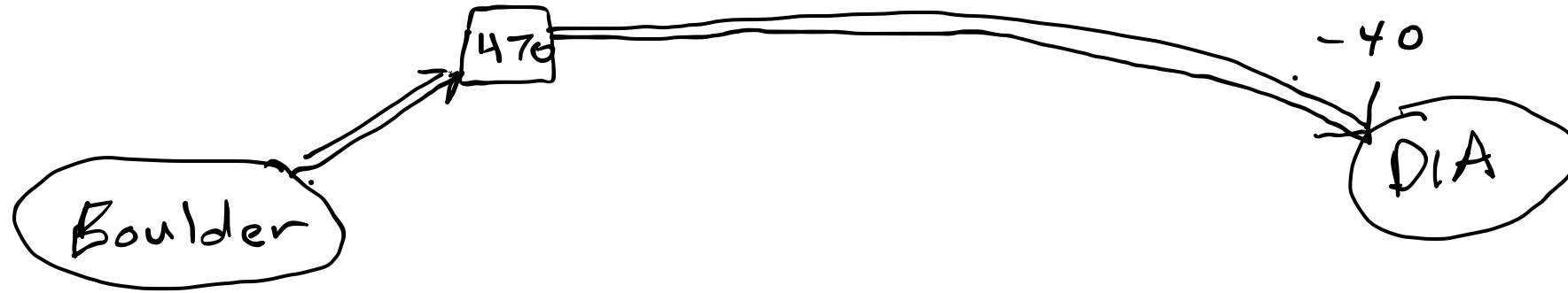


Breakout Rooms: DIA Run

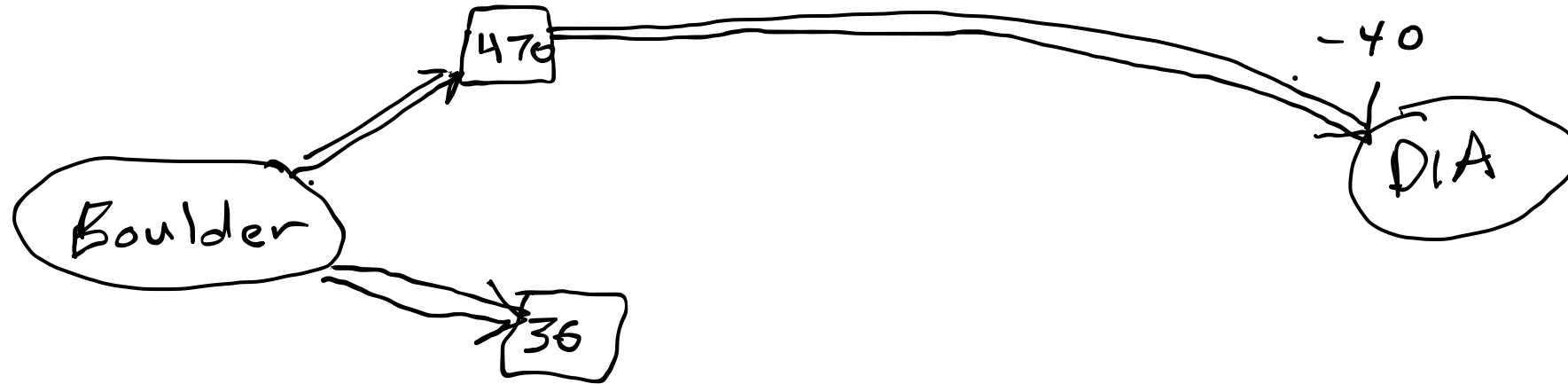
Boulder

DIA

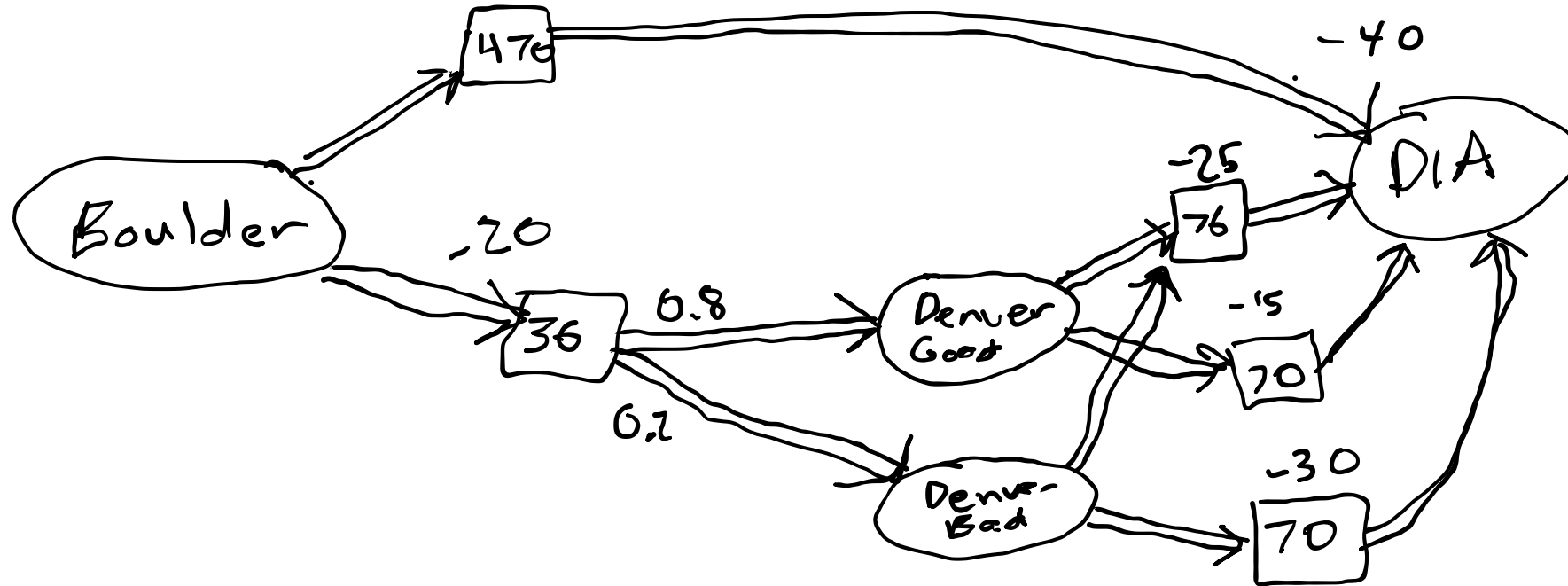
Breakout Rooms: DIA Run



Breakout Rooms: DIA Run



Breakout Rooms: DIA Run



Value Iteration

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

initialize V, V' (differently)

while $\|V - V'\|_\infty < \epsilon$

$V \leftarrow V'$

$V'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) V^\pi(s')) \quad \forall s \in S$

return $V' \approx V^*$

$$\|V\|_\infty = \max_i V_i$$

apply Bellman's

$$\|V' - V^*\|_\infty < \epsilon$$

HW2

Guiding Questions

Guiding Questions

$$V^*(s)$$

$$V(s)$$

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

Value Iter — apply Bellman's
Policy Iter < eval
update