

**Last Time**

# Last Time

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

$\phi$  initial state dist.  
 $(S, A, T, R, \gamma)$

# Last Time

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

# Last Time

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

# 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**?

# Value-Based Policy Evaluation

Want  $U(\pi)$

$$\begin{aligned}
 U^\pi(s) &= E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s \right] \\
 &= E[R(s, \pi(s))] + E \left[ \sum_{t=1}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_1 = s' \right] \\
 &\quad s' \sim T(s, \pi(s)) \\
 &= R(s, \pi(s)) + \gamma E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_1 = s' \right] \\
 &= R(s, \pi(s)) + \gamma E[U^\pi(s')] \\
 &\quad s' \sim T(s, \pi(s))
 \end{aligned}$$

Bellman Expectation Equation

$$\begin{aligned}
 &\rightarrow \boxed{U^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s' | s, \pi(s)) U^\pi(s')} \\
 &\quad \bar{U}^\pi = \bar{R}^\pi + \gamma \bar{T}^\pi \bar{U}^\pi
 \end{aligned}$$

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

$$\boxed{\bar{U}^\pi = (I - \gamma \bar{T}^\pi)^{-1} \bar{R}^\pi}$$

$$\bar{U}_i^\pi = U^\pi(i)$$

$$\bar{T}_{ij}^\pi = T(j | i, \pi(i))$$

$$\bar{R}_i^\pi = R(i, \pi(i))$$

$$\begin{aligned}
 U(\pi) &= \sum_s b(s) U^\pi(s) \\
 &\quad \bar{U}^\pi \cdot \bar{b}
 \end{aligned}$$

$$\bar{b}_i = b(s)$$

# MDP Example: Up-Down Problem

For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)

1



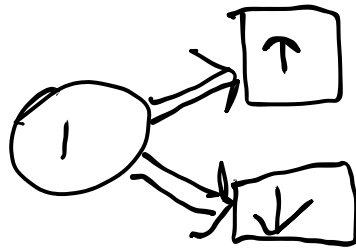
# MDP Example: Up-Down Problem

For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)



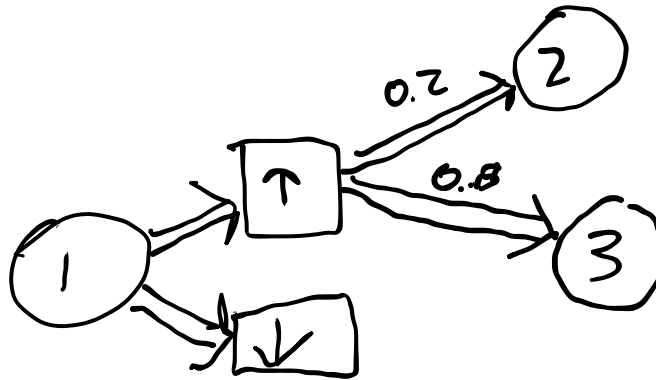
# MDP Example: Up-Down Problem

For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)



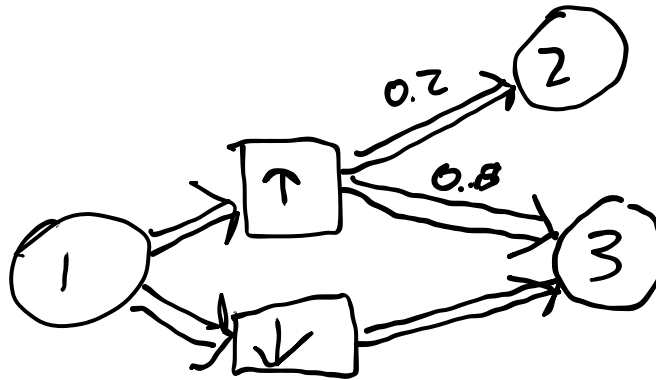
# MDP Example: Up-Down Problem

For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)



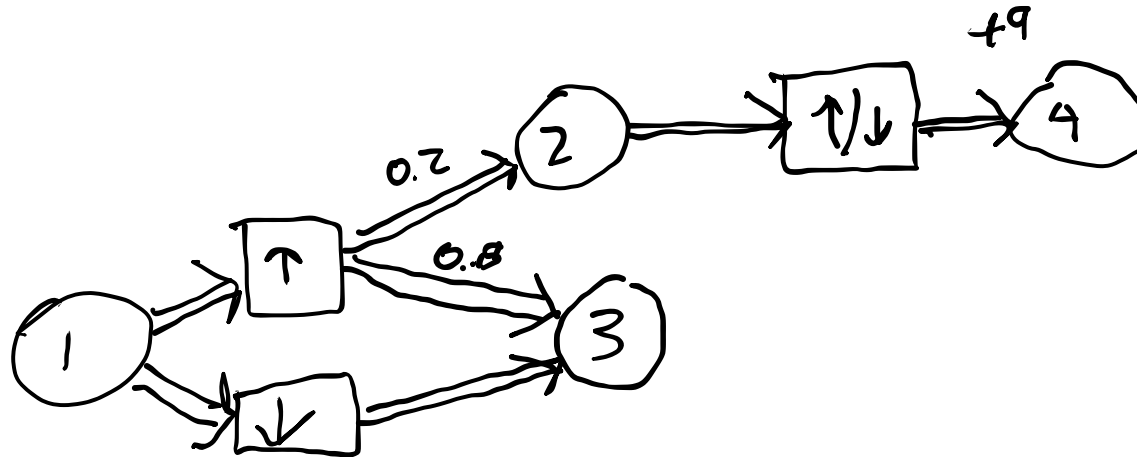
# MDP Example: Up-Down Problem

For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)



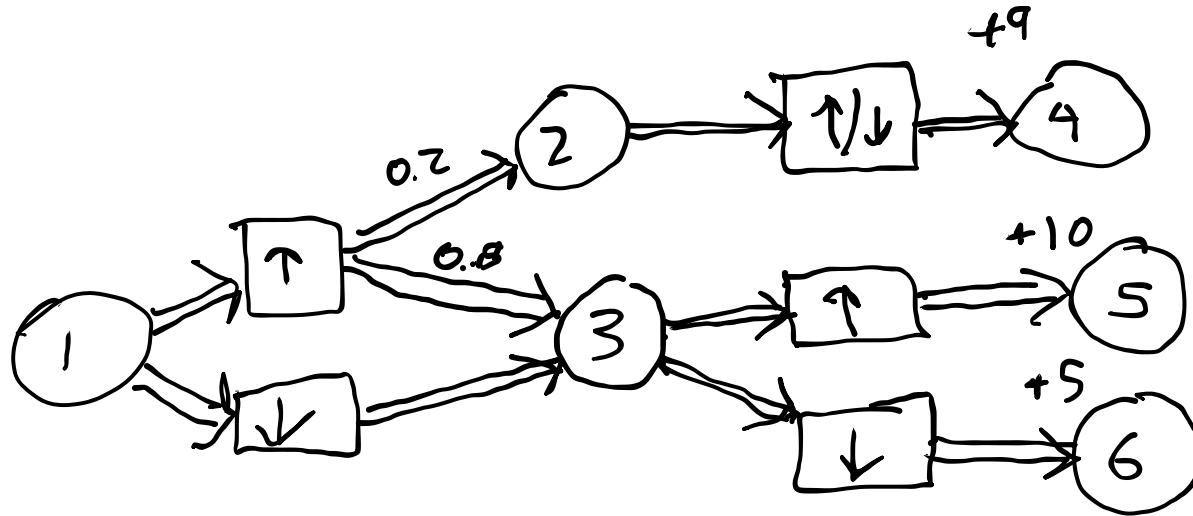
# MDP Example: Up-Down Problem

For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)

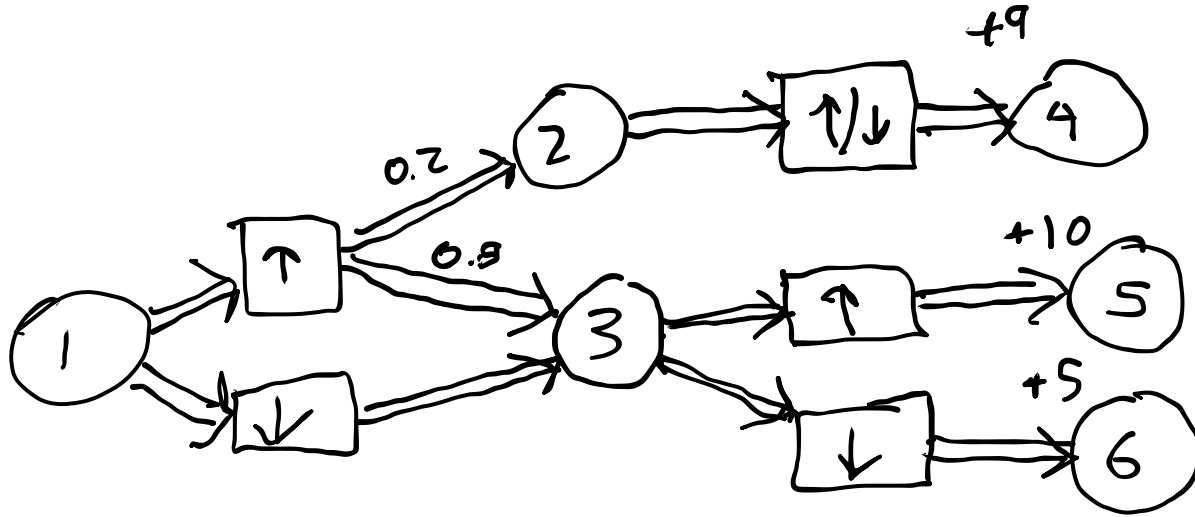


# MDP Example: Up-Down Problem

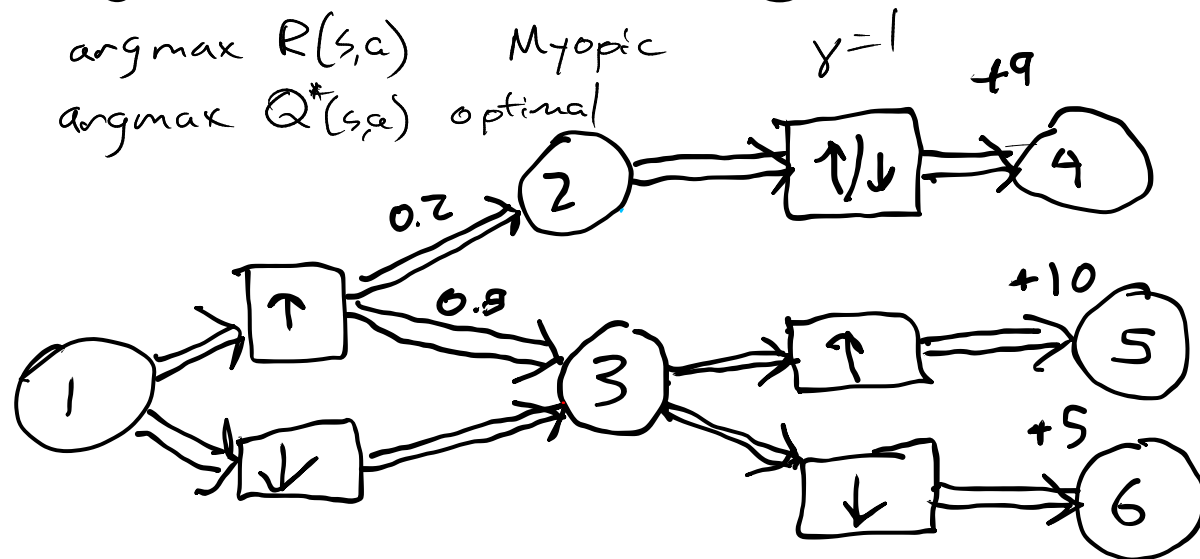
For this lecture,  $\Rightarrow$  is same as  $\rightarrow$  (distinguishes from Bayes Net)



# Dynamic Programming and Value Backup



# Dynamic Programming and Value Backup



Bellman Backup

$U^*(s) \leftarrow 0$  for all terminal states

Repeat until all  $U^*(s)$  calculated

find  $\pi^*, U^*$  for states where  $U^*$  is known for all successor states

Only works if there are no cycles?

Bellman's Principle of Optimality: Every sub-policy in an optimal policy is locally optimal

$$U^*(s) = \max_{\pi} U^{\pi}(s) \quad \leftarrow \text{optimal value function}$$

$$\pi^+(s) = \operatorname{argmax}_{a \in A} \underbrace{(R(s,a) + \gamma E[U^*(s')])}_{Q^*(s,a)}$$

s	$U^*(s)$	a	$Q^*(s,a)$
4	0		
5	0		
6	0		
2	9	U/D	$+9 + \gamma \cdot 0$
3	10	U	$+10 + \gamma \cdot 0$
		D	$+5 + \gamma \cdot 0$
1	10	U	$0 + \gamma(0.2 U^*(2) + 0.8 U^*(3))$
		D	$0 + \gamma(U^*(3) + 10)$



# Break: DIA Run

Boulder.

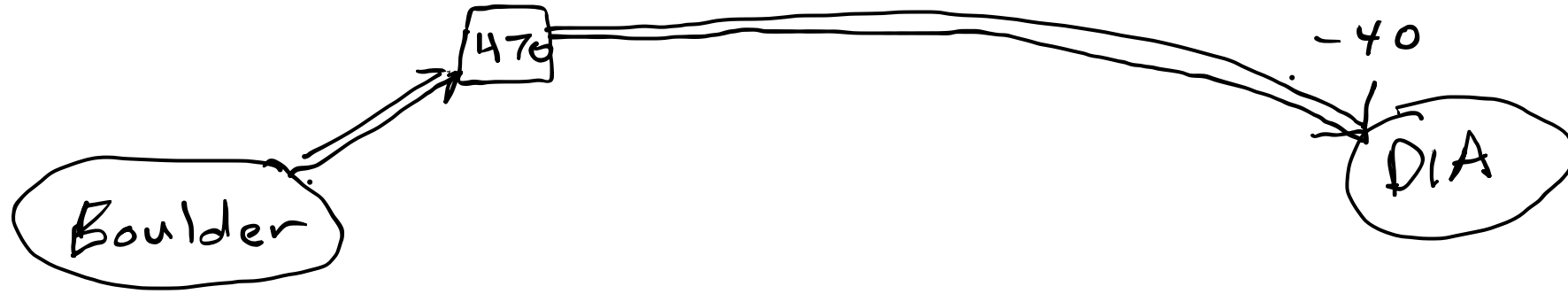


# Break: DIA Run

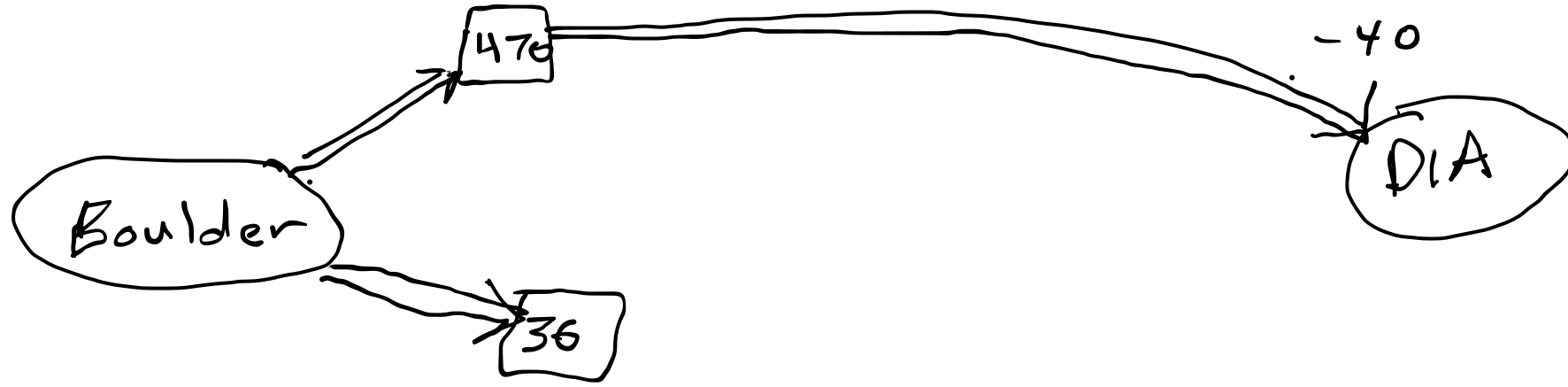
Boulder

DIA

# Break: DIA Run



# Break: DIA Run

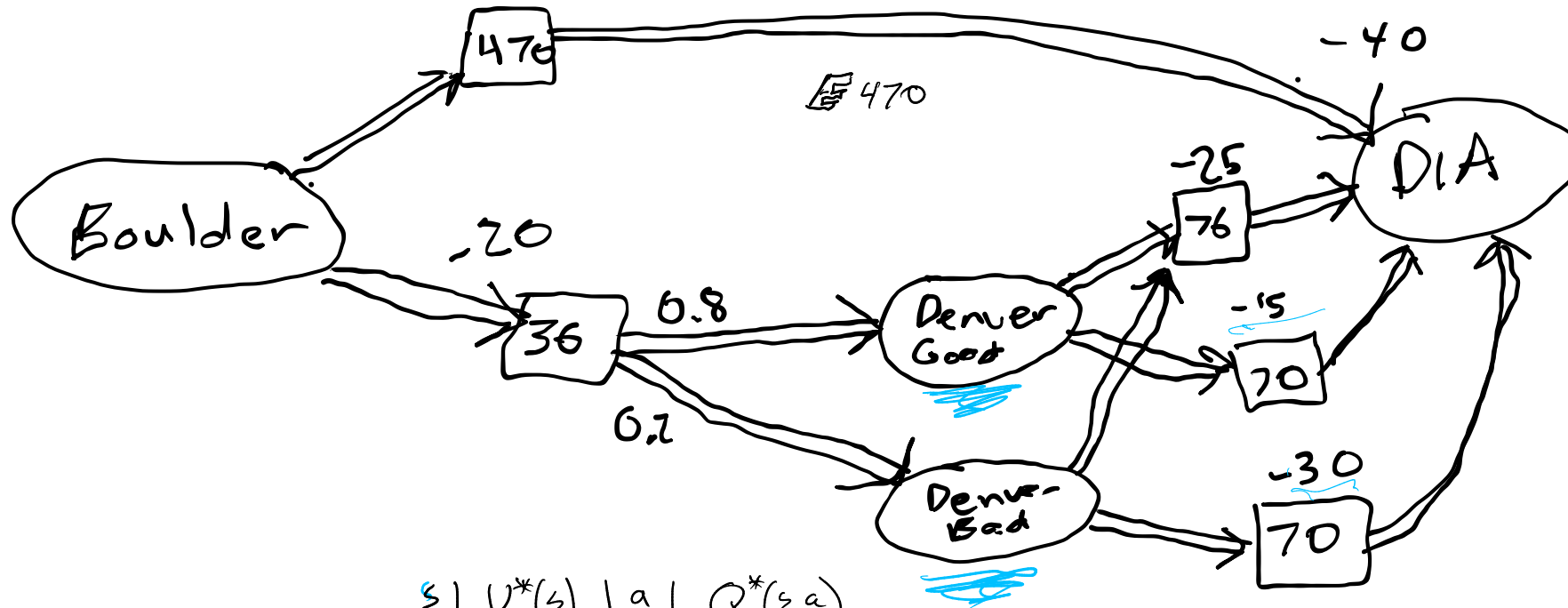


# Break: DIA Run

Myopic

$$\operatorname{argmax} R(\text{Boulder}, a)$$

$$= \boxed{36}$$



$s$	$U^*(s)$	$a$	$Q^*(s, a)$
DIA	0		
Den Good	-15	76 <u>70</u>	-25 -15
Den Bad	-25	76 70	-25 -30
Boulder	-37	470 <u>36</u>	-40 $-20 + 0.2(-25) + 0.8(-15) = -37$

# Policy Iteration

Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

# Policy Iteration

Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

1. initialize  $\pi, \pi'$  (differently)

# Policy Iteration

## Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

1. initialize  $\pi, \pi'$  (differently)
2. while  $\pi \neq \pi'$



# Policy Iteration

## Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

1. initialize  $\pi, \pi'$  (differently)
2. while  $\pi \neq \pi'$
3.    $\pi \leftarrow \pi'$

# Policy Iteration

$U^*$

## Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

1. initialize  $\pi, \pi'$  (differently)
2. while  $\pi \neq \pi'$
3.  $\pi \leftarrow \pi'$
4.  $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$

$$\pi'(s) = \underset{a}{\operatorname{argmax}} \left( R(s,a) + \gamma E[U^\pi(s)] \right)$$

# Policy Iteration

## Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

1. initialize  $\pi, \pi'$  (differently)
2. while  $\pi \neq \pi'$
3.    $\pi \leftarrow \pi'$
4.    $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$
5.    $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$

# Policy Iteration

## Algorithm: Policy Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$

1. initialize  $\pi, \pi'$  (differently)

2. while  $\pi \neq \pi'$

3.  $\pi \leftarrow \pi'$

4.  $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$

5.  $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$

6. return  $\pi$

swap  
for MC  
Evaluation

1. Policy Evaluation

2. Policy Improvement

# Value Iteration

Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

# Value Iteration

Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)

# Value Iteration

## Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)
2. while  $\|U - U'\|_{\infty} < \epsilon$

# Value Iteration

## Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)
2. while  $\|U - U'\|_{\infty} < \epsilon$
3.  $U \leftarrow U'$



# Value Iteration

## Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)
2. while  $\|U - U'\|_\infty < \epsilon$
3.  $U \leftarrow U'$
4.  $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\bullet(s')) \quad \forall s \in S$

# Value Iteration

## Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)
2. while  $\|U - U'\|_\infty < \epsilon$
3.  $U \leftarrow U'$
4.  $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U^\pi(s')) \quad \forall s \in S$
5. return  $U'$

# Value Iteration

## Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)
2. while  $\|U - U'\|_\infty < \epsilon$
3.  $U \leftarrow U'$
4.  $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$
5. return  $U'$

- Returned  $U'$  will be <sup>near</sup>  $U^*$ !

# Value Iteration

## Algorithm: Value Iteration

Given: MDP  $(S, A, R, T, \gamma, b)$ , tolerance  $\epsilon$

1. initialize  $U, U'$  (differently)
2. while  $\|U - U'\|_\infty < \epsilon$
3.  $U \leftarrow U'$
4.  $U'(s) \leftarrow \max_{a \in A} \left( \overset{\text{Immediate Reward}}{R(s, a)} + \gamma \sum_{s' \in S} \overset{\text{Future Value}}{T(s'|s, a)U^\pi(s')} \right) \quad \forall s \in S$
5. return  $U'$

- Returned  $U'$  will be  $U^*$ !
- $\pi^*$  is easy to extract:  $\pi^*(s) = \arg \max (R(s, a) + \gamma E[U^*(s)])$

# Bellman's Equations

Policy Evaluation  
(Linear)

Bellman Backup  
Certificate of  
Optimality

Value  
Iteration

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

Bellman's Expectation  
Equation

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

Bellman's Optimality "Bellman's  
Equation"

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

Bellman  
Operator

VI:  
initialize  $U, U'$   
while  $\|U - U'\| > \epsilon$   
     $U \leftarrow U'$   
     $U' \leftarrow B[U]$

# 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**?

$$Q(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim \pi(s,a)} [V(s')]$$

Value Iteration = Repeatedly Applying Bellman Operation

*Handwritten annotations:*  $R(s,a)$  is underlined in blue and labeled "immediate";  $\mathbb{E}_{s' \sim \pi(s,a)} [V(s')]$  is underlined in red and labeled "future".

Policy Iteration = Repeatedly  
Evaluating Policy  
Improving Policy