

- 
- Going forward
  - Gridworld
  - Q1: Not in-place

# Utilities of Sequences

---

- In order to formalize optimality of a policy, need to understand utilities of sequences of rewards
- Typically consider **stationary preferences**:

$$\begin{aligned} [r, r_0, r_1, r_2, \dots] \succ [r, r'_0, r'_1, r'_2, \dots] \\ \Leftrightarrow \\ [r_0, r_1, r_2, \dots] \succ [r'_0, r'_1, r'_2, \dots] \end{aligned}$$

- **Theorem: only two ways to define stationary utilities**
  - Additive utility:

$$U([r_0, r_1, r_2, \dots]) = r_0 + r_1 + r_2 + \dots$$

- Discounted utility:

$$U([r_0, r_1, r_2, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 \dots$$

# Infinite Utilities?!

- Problem: infinite state sequences have infinite rewards

- Solutions:

- Finite horizon:

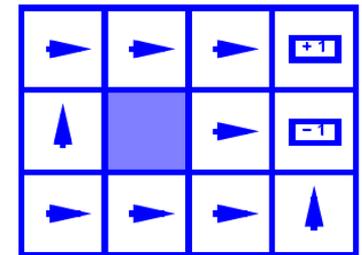
- Terminate episodes after a fixed T steps (e.g. life)
- Gives nonstationary policies ( $\pi$  depends on time left)

- Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like “done” for High-Low)

- Discounting: for  $0 < \gamma < 1$

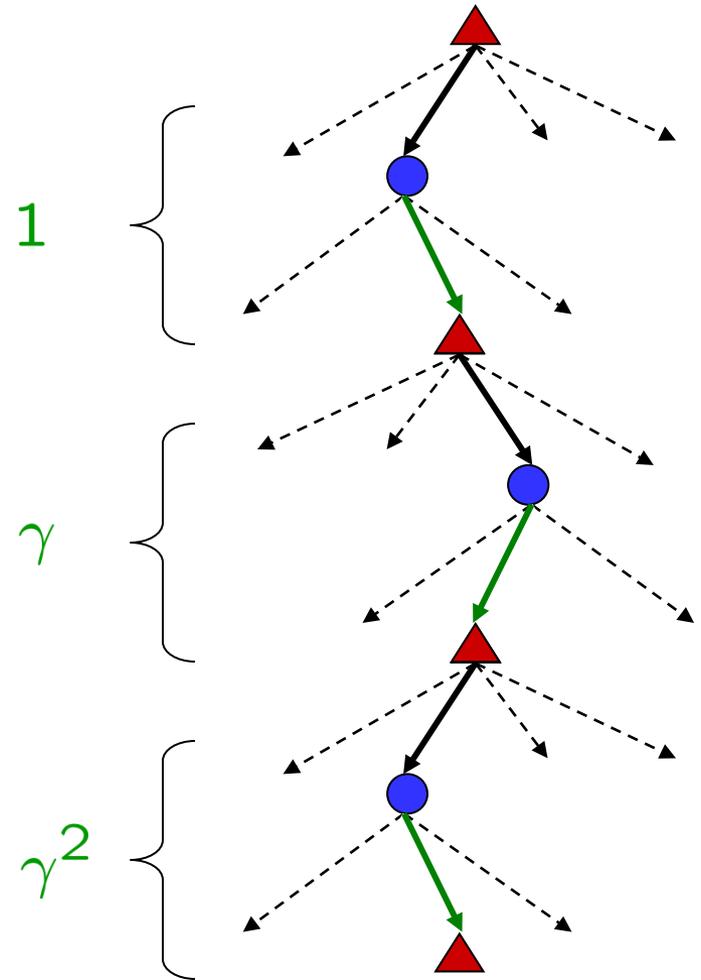
$$U([r_0, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max} / (1 - \gamma)$$

- Smaller  $\gamma$  means smaller “horizon” – shorter term focus



# Discounting

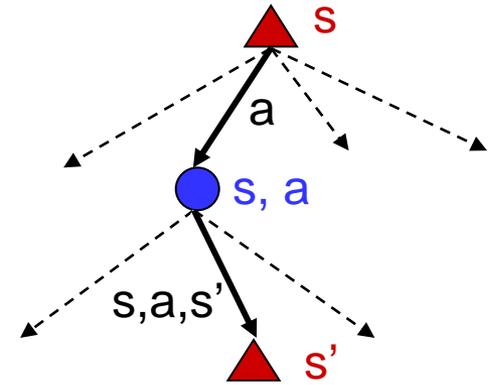
- Typically discount rewards by  $\gamma < 1$  each time step
  - Sooner rewards have higher utility than later rewards
  - Also helps the algorithms converge



# Recap: Defining MDPs

---

- Markov decision processes:
  - States  $S$
  - Start state  $s_0$
  - Actions  $A$
  - Transitions  $P(s'|s,a)$  (or  $T(s,a,s')$ )
  - Rewards  $R(s,a,s')$  (and discount  $\gamma$ )



- MDP quantities so far:
  - Policy = Choice of action for each state
  - Utility (or return) = sum of discounted rewards

# Optimal Utilities

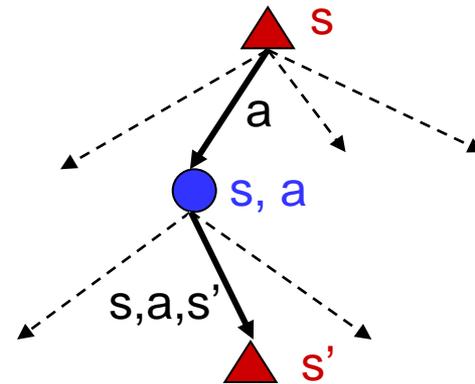
- Fundamental operation: compute the values (optimal expectimax utilities) of states  $s$

- Why? Optimal values define optimal policies!

- Define the value of a state  $s$ :  
 $V^*(s)$  = expected utility starting in  $s$  and acting optimally

- Define the value of a q-state  $(s,a)$ :  
 $Q^*(s,a)$  = expected utility starting in  $s$ , taking action  $a$  and thereafter acting optimally

- Define the optimal policy:  
 $\pi^*(s)$  = optimal action from state  $s$



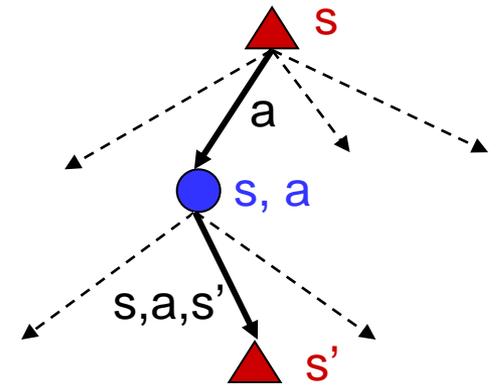
3	0.812	0.868	0.912	<span style="border: 1px solid black; padding: 2px;">+1</span>
2	0.762		0.660	<span style="border: 1px solid black; padding: 2px;">-1</span>
1	0.705	0.655	0.611	0.388
	1	2	3	4

3	→	→	→	<span style="border: 1px solid black; padding: 2px;">+1</span>
2	↑		↑	<span style="border: 1px solid black; padding: 2px;">-1</span>
1	↑	←	←	←
	1	2	3	4

# The Bellman Equations

- Definition of “optimal utility” leads to a simple one-step lookahead relationship amongst optimal utility values:

Optimal rewards = maximize over first action and then follow optimal policy



- Formally:

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

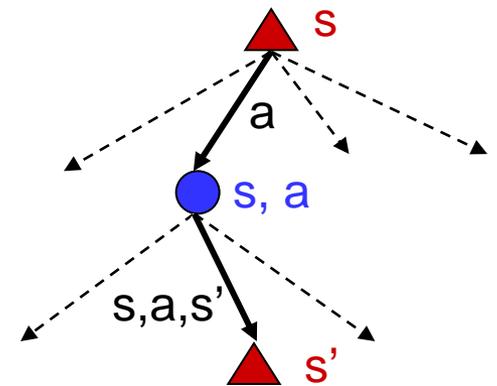
# Solving MDPs

- We want to find the **optimal policy**  $\pi^*$
- Proposal 1: modified expectimax search, starting from each state  $s$ :

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

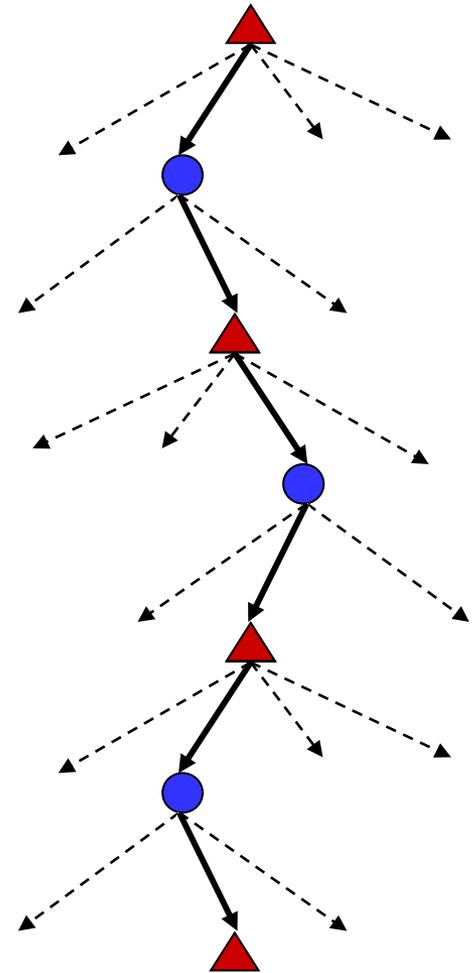
$$V^*(s) = \max_a Q^*(s, a)$$



# Why Not Search Trees?

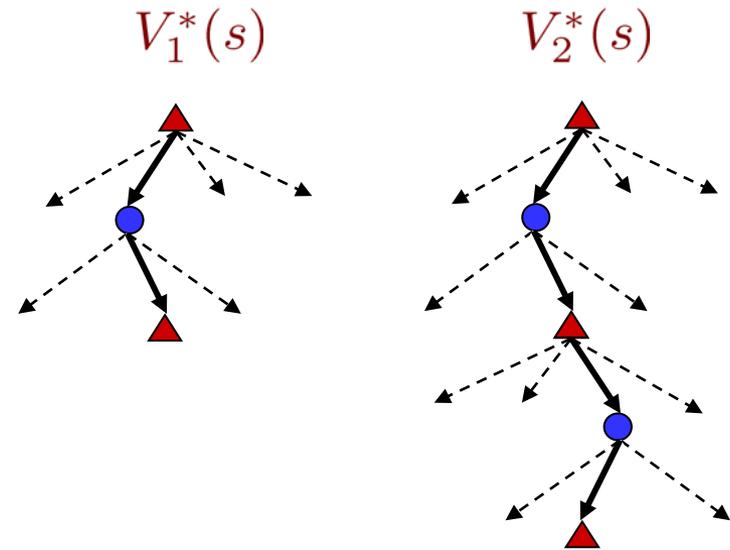
---

- Why not solve with expectimax?
- Problems:
  - This tree is usually infinite (why?)
  - Same states appear over and over (why?)
  - We would search once per state (why?)
- Idea: Value iteration
  - Compute optimal values for all states all at once using successive approximations
  - Will be a bottom-up dynamic program similar in cost to memoization
  - Do all planning offline, no replanning needed!



# Value Estimates

- Calculate estimates  $V_k^*(s)$ 
  - Not the optimal value of  $s$ !
  - The optimal value considering only next  $k$  time steps ( $k$  rewards)
  - As  $k \rightarrow \infty$ , it approaches the optimal value
- Almost solution: recursion (i.e. expectimax)
- Correct solution: dynamic programming



# Value Iteration

---

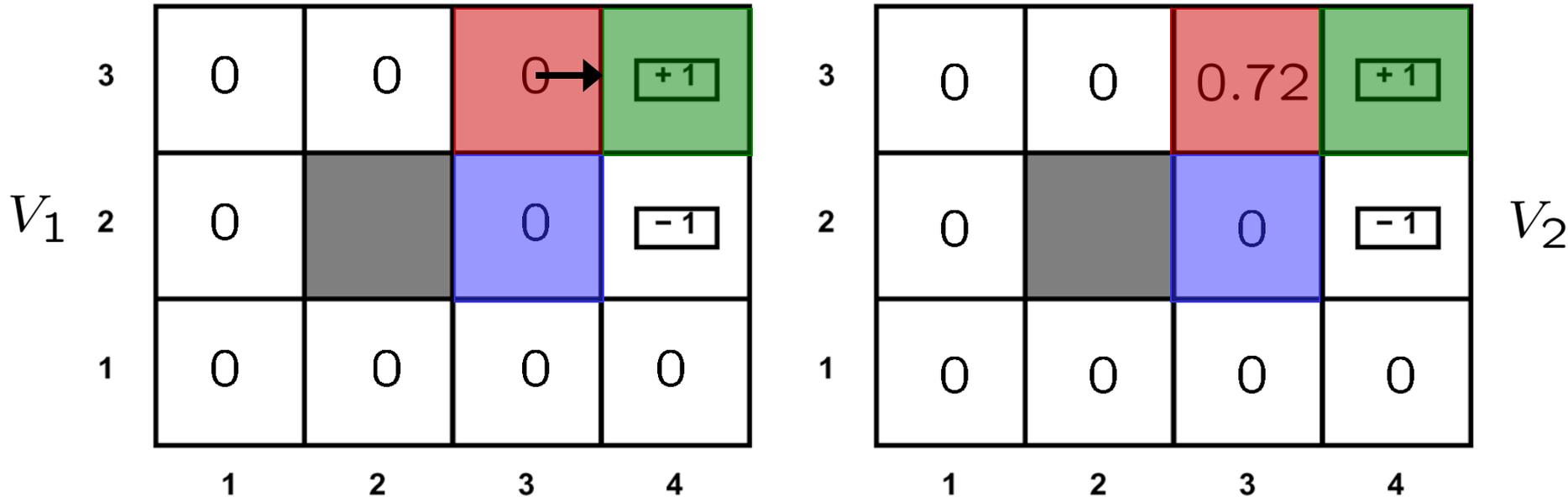
- **Idea:**

- Start with  $V_0^*(s) = 0$ , which we know is right (why?)
- Given  $V_i^*$ , calculate the values for all states for depth  $i+1$ :

$$V_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

- This is called a **value update** or **Bellman update**
  - Repeat until convergence
- **Theorem: will converge to unique optimal values**
    - Basic idea: approximations get refined towards optimal values
    - Policy may converge long before values do

# Example: Bellman Updates



$$V_{i+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

$$V_2(\langle 3, 3 \rangle) = \sum_{s'} T(\langle 3, 3 \rangle, \text{right}, s') [R(\langle 3, 3 \rangle) + 0.9 V_1(s')]$$

max happens for  
 $a=\text{right}$ , other  
actions not shown

$$= 0.9 [0.8 \cdot 1 + 0.1 \cdot 0 + 0.1 \cdot 0]$$

# Example: Value Iteration

---

$V_2$

3	0	0	0.72	<span style="border: 1px solid black; padding: 2px;">+1</span>
2	0		0	<span style="border: 1px solid black; padding: 2px;">-1</span>
1	0	0	0	0
	1	2	3	4

$V_3$

3	0	0.52	0.78	<span style="border: 1px solid black; padding: 2px;">+1</span>
2	0		0.43	<span style="border: 1px solid black; padding: 2px;">-1</span>
1	0	0	0	0
	1	2	3	4

- Information propagates outward from terminal states and eventually all states have correct value estimates

# Convergence\*

---

- Define the max-norm:  $\|U\| = \max_s |U(s)|$
- Theorem: For any two approximations U and V

$$\|U^{t+1} - V^{t+1}\| \leq \gamma \|U^t - V^t\|$$

- I.e. any distinct approximations must get closer to each other, so, in particular, any approximation must get closer to the true U and value iteration converges to a unique, stable, optimal solution
- Theorem:

$$\|U^{t+1} - U^t\| < \epsilon, \Rightarrow \|U^{t+1} - U\| < 2\epsilon\gamma/(1 - \gamma)$$

- I.e. once the change in our approximation is small, it must also be close to correct

# Practice: Computing Actions

---

- Which action should we chose from state  $s$ :
  - Given optimal values  $V$ ?

$$\arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

- Given optimal q-values  $Q$ ?

$$\arg \max_a Q^*(s, a)$$

- Lesson: actions are easier to select from  $Q$ 's!

# Utilities for Fixed Policies

- Another basic operation: compute the utility of a state  $s$  under a fixed (general non-optimal) policy
- Define the utility of a state  $s$ , under a fixed policy  $\pi$ :
  - $V^\pi(s)$  = expected total discounted rewards (return) starting in  $s$  and following  $\pi$
- Recursive relation (one-step look-ahead / Bellman equation):

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

