

6.390 Intro to Machine Learning

Lecture 11: Markov Decision Processes

Shen Shen

Nov 13, 2025

11am, Room 10-250

Interactive Slides and Lecture Recording

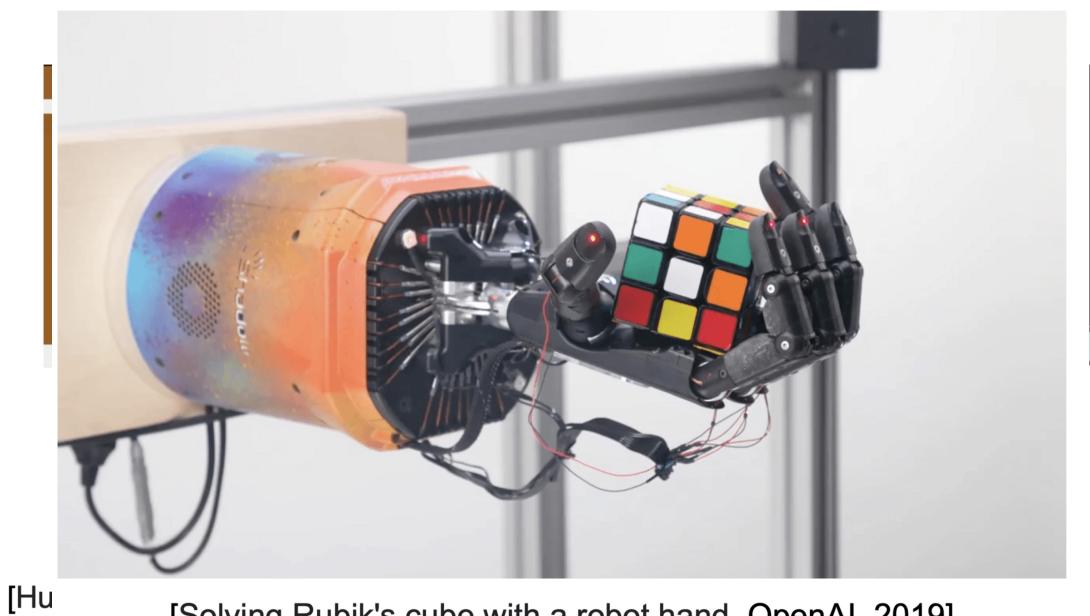
1

Learning to Walk

Massachusetts Institute of Technology, 2004



Toddler demo, Russ Tedrake thesis, 2004 uses vanilla policy gradient (actor-critic)



[Solving Rubik's cube with a robot hand. OpenAl. 2019]

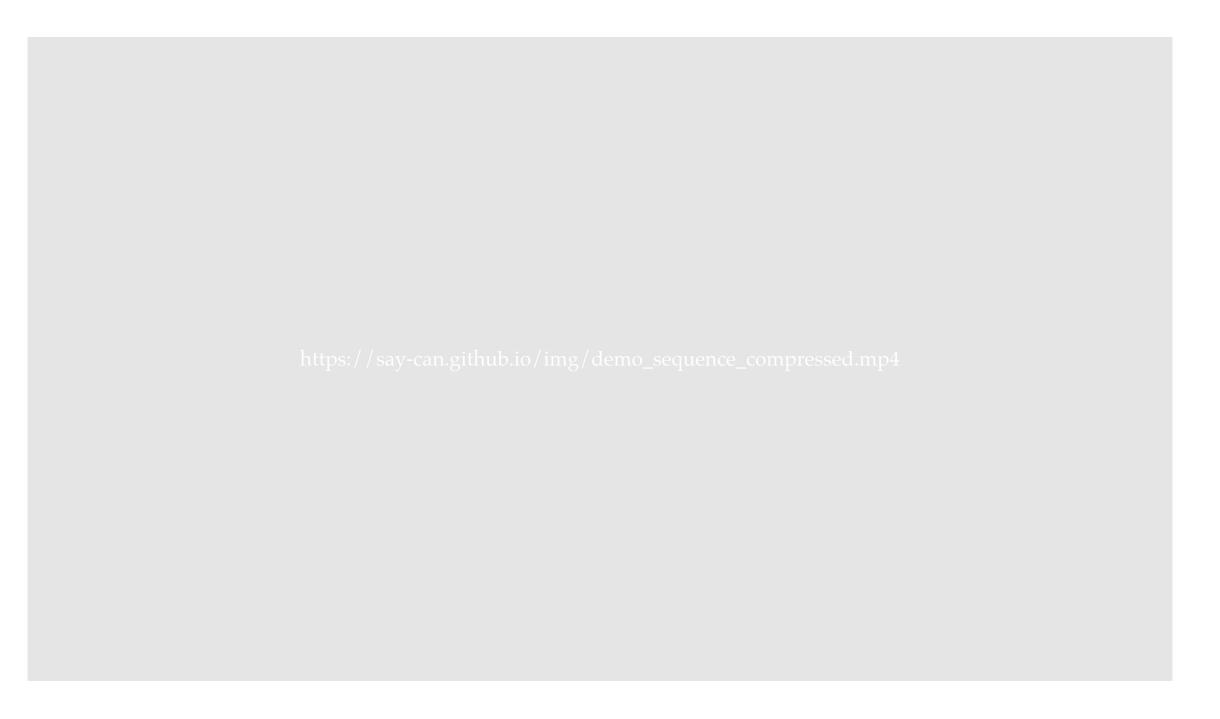
15]

Discovering faster that rix multiplication

algo ChatGPT: Optimizing Language Models for Dialogue https://doi.org/ Received: 2 Oc BAccepted: 2 Au Published onlir We've trained a model called ChatGPT which interact. a CONVERSATIONAL WAY. The dialogue format makes it possesses to the minimum of the ChatGPT to answer followup questions, admit its mistake shallenge incorrect premises, and reject inappro G. ChatGPT is a sibling model to InstructGPT a_1

YEDEI

[Aligning language models to follow instructions. Ouyang et al. 2022]



Outline

- Markov Decision Processes Definition
- Policy Evaluation
 - State value functions: V^{π}
 - Bellman recursions and Bellman equations
- Policy Optimization
 - Optimal policies π^*
 - Optimal action value functions: Q*
 - Value iteration

Markov Decision Processes

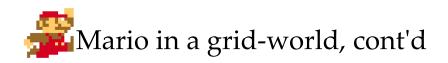
- Research area initiated in the 50s by Bellman, known under various names:
 - Stochastic optimal control (Control theory)
 - Stochastic shortest path (Operations research)
 - Sequential decision making under uncertainty (Economics)
 - Reinforcement learning (Artificial intelligence, Machine learning)
- A rich variety of elegant theory, mathematics, algorithms, and applications—but also considerable variation in notation.
- We will use the most RL-flavored notations.



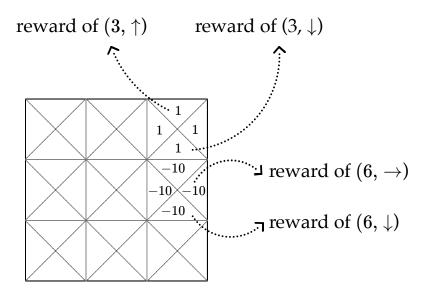
1	2	3 ↑ 80%
4	20% ··· 5	6
7	8	9

Running example: Mario in a grid-world

- 9 possible states s
- 4 possible actions a: {Up \uparrow , Down \downarrow , Left \leftarrow , Right \rightarrow }
- (state, action) results in a **transition** T into a next state:
 - Normally, we get to the "intended" state;
 - E.g., in state (7), action "↑" gets to state (4)
 - If an action would take Mario out of the grid world, stay put;
 - \circ E.g., in state (9), " \rightarrow " gets back to state (9)
 - In state (6), action "↑" leads to two possibilities:
 - 20% chance to (2)
 - 80% chance to (3).



- (state, action) pairs give **rewards**:
- **⚠** in state 3, any action gives reward 1
- ♠ in state 6, any action gives reward -10
- any other (state, action) pair gives reward 0



- **discount factor**: a scalar of 0.9 that reduces the 'worth' of future rewards depending on when Mario receives them.
 - So, e.g., for $(3, \leftarrow)$ pair, Mario gets
 - o at the start of the game, a reward of 1
 - at the 2nd time step, a discounted reward of 0.9
 - \circ at the 3rd time step, a further discounted reward of $(0.9)^2$... and so on

- S: state space, contains all possible states s.
- \mathcal{A} : action space, contains all possible actions a.

In 6.390,

• \mathcal{S} and \mathcal{A} are small discrete sets, unless otherwise specified.

- S: state space, contains all possible states s.
- \mathcal{A} : action space, contains all possible actions a.
- T (s, a, s'): the probability of transition from state s to s' when action a is taken.

1	2 ▼	3 80%
4	5	6
7	8	9

$$\mathrm{T}\left(7,\uparrow,4
ight)=1$$

$$T(9, \rightarrow, 9) = 1$$

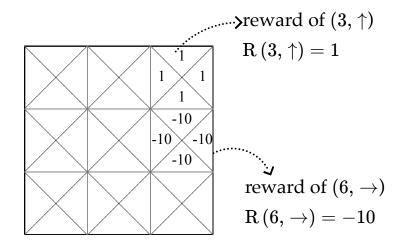
$$T(6,\uparrow,3) = 0.8$$

$$\mathrm{T}\left(6,\uparrow,2
ight)=0.2$$

In 6.390,

- \mathcal{S} and \mathcal{A} are small discrete sets, unless otherwise specified.
- *s'* and *a'* are short-hand for the next-timestep state and action.

- S: state space, contains all possible states s.
- \mathcal{A} : action space, contains all possible actions a.
- T(s, a, s'): the probability of transition from state s to s' when action a is taken.
- R(s, a): reward, takes in a (state, action) pair and returns a reward.



In 6.390,

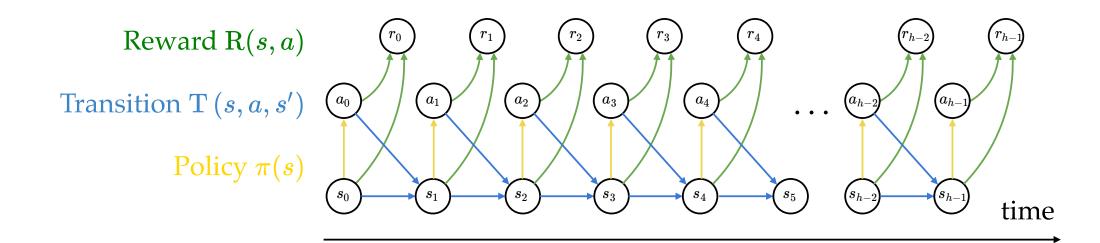
- \mathcal{S} and \mathcal{A} are small discrete sets, unless otherwise specified.
- *s'* and *a'* are short-hand for the next-timestep state and action.
- R(s, a) is deterministic and bounded.

- S: state space, contains all possible states s.
- \mathcal{A} : action space, contains all possible actions a.
- T (s, a, s'): the probability of transition from state s to s' when action a is taken.
- R(s, a): reward, takes in a (state, action) pair and returns a reward.
- $\gamma \in [0,1]$: discount factor, a scalar.

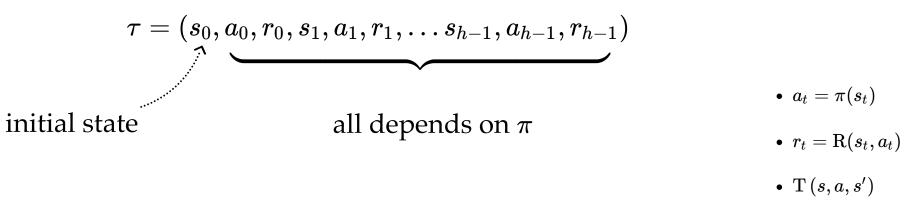
• $\pi(s)$: policy, takes in a state and returns an action. The goal of an MDP is to find a good policy.

In 6.390,

- \mathcal{S} and \mathcal{A} are small discrete sets, unless otherwise specified.
- *s'* and *a'* are short-hand for the next-timestep state and action.
- R(s, a) is deterministic and bounded.
- $\pi(s)$ is deterministic.



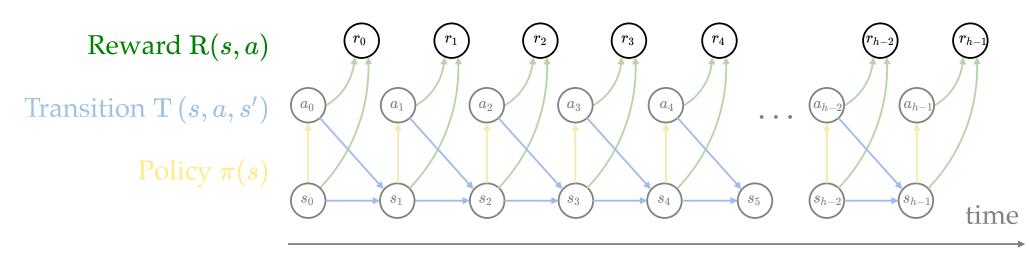
a trajectory (also called an experience or rollout) of horizon h



Outline

- Markov Decision Processes Definition
- Policy Evaluation
 - State value functions: V^{π}
 - Bellman recursions and Bellman equations
- Policy Optimization
 - Optimal policies π^*
 - Optimal action value functions: Q*
 - Value iteration

Starting in a given s_0 , how good is it to follow a *given* policy π for h time steps?



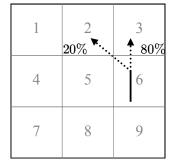
One idea:

$$\mathrm{R}(s_0,\pi(s_0)) \,+\, \gamma \mathrm{R}(s_1,\pi(s_1)) \,+\, \gamma^2 \mathrm{R}(s_2,\pi(s_2)) \,+\, \gamma^3 \mathrm{R}(s_3,\pi(s_3)) \,+\, \ldots \,+\, \gamma^{h-1} \mathrm{R}(s_{h-1},\pi(s_{h-1}))$$

But if we start at $s_0 = 6$ and follow the "always-up" policy:

states and one special transition:

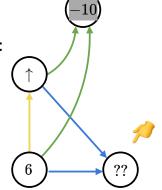




rewards:

	0			0		1 /	/
0	\times	0	0	\times	0	1 🗙	1
	0		/	0		/1\	\
	0	/		0	/	-10/	/
0	\times	0	0	X	0	-10×-1	0
	0		/	0		/-10	\
	0	/	$\overline{}$	0	/	0	/
0	\times	0	0	\times	0	0×0)
	0		_	0		/ 0 \	\
	0						$\begin{array}{cccccccccccccccccccccccccccccccccccc$

trajectory:



Starting in a given s_0 , how good is it to follow a *given* policy π for h time steps?

Value functions:

$$\begin{aligned} \mathbf{V}_h^{\pi}(s) &:= \mathbb{E} \Big[\mathbf{R}(s_0, \pi(s_0)) \, + \, \gamma \mathbf{R}(s_1, \pi(s_1)) \, + \, \gamma^2 \mathbf{R}(s_2, \pi(s_2)) \, + \, \gamma^3 \mathbf{R}(s_3, \pi(s_3)) \, + \, \dots \, + \, \gamma^{h-1} \mathbf{R}(s_{h-1}, \pi(s_{h-1})) \Big] \\ &= \mathbb{E} \left[\sum_{t=0}^{h-1} \gamma^t \mathbf{R}\left(s_t, \pi\left(s_t\right)\right) \, | \, s_0 = s, \pi \right] \end{aligned} \quad \text{(eq. 1)}$$

$$\text{in 6.390, this expectation is only w.r.t. the transition probabilities T} \left(s, a, s'\right)$$

- $V_h^{\pi}(s)$: expected sum of discounted rewards starting in state s and follow π for h steps
- Value is long-term; reward is immediate (one-time)
- Horizon-0 values $V_0^{\pi}(s)$ defined as 0 for all states



evaluate $V_h^\pi(s)$ under the "always-up" policy

states and one special transition:

1	2 20% ···.	3 ≜ 80%
4	5	6
7	8	9

rewards

0 /	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	1
0 /	0 /	-10/
0×0	0×0	-10 <-10
0	0	-10
0 /	0 /	0 /
0×0	0×0	$ 0 \times 0 $
0	0	0

- $\pi(s) = ``\uparrow", \ \forall s$
- $\gamma = 0.9$

$\mathrm{V}_h^{\uparrow}(s) = \mathbb{E}\left[\sum_{t=0}^{h-1} \gamma^t \mathrm{R}\left(s_t, \uparrow ight) \mid s_0 = s ight]$
$=\mathbb{E}\Big[\mathrm{R}(s_0,\uparrow)+\gamma\mathrm{R}(s_1,\uparrow)+\cdots+\gamma^{h-1}\mathrm{R}(s_{h-1},\uparrow)\Big]$
h terms

horizon h = 0: no step left

horizon h = 1: receive the rewards



horizon h = 2

states and

one special transition:

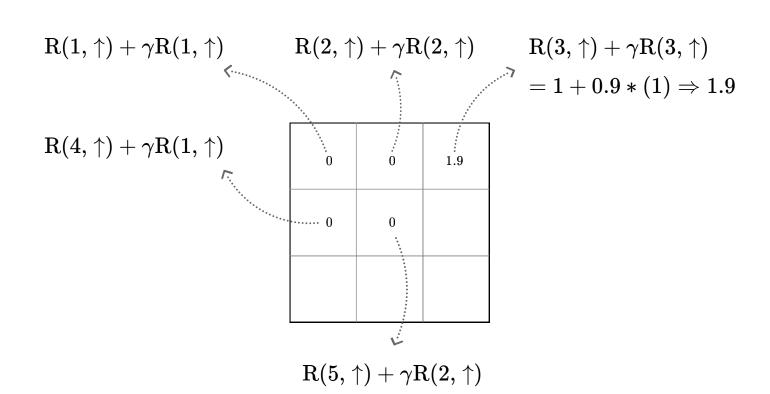
1	2 20% ···	3 ≜ 80%
4	5	6
7	8	9

0 /	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	/1
0 /	0 /	-10/
0×0	0×0	-10 \(\sigma -10 \)
0	0	-10
0 /	0 /	0 /
0×0	0×0	$ 0 \times 0 $
0	0	0

•
$$\pi(s) = ``\uparrow", \forall s$$

•
$$\gamma=0.9$$

$$\mathrm{V}_2^{\uparrow}(s): \mathbb{E}\Big[\underbrace{\mathrm{R}(s_0,\uparrow) + \gamma\mathrm{R}(s_1,\uparrow)}_{2 ext{ terms}}\Big]$$





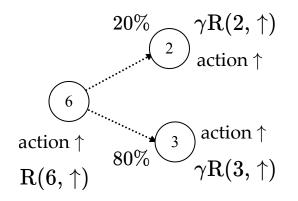
horizon h = 2

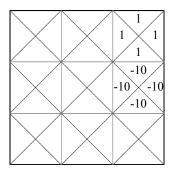
states and

one special transition:

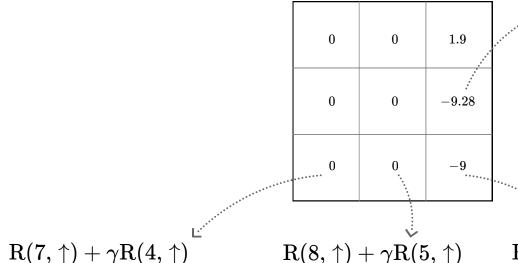
1	2 20% [▼] ···.	3 • <u>*</u> 80%
4	5	6
7	8	9

$$\mathrm{V}_2^{\uparrow}(s): \mathbb{E}\Big[\underbrace{\mathrm{R}(s_0,\uparrow) + \gamma \mathrm{R}(s_1, \uparrow)}_{2 ext{ terms}} \Big]$$





- $\pi(s) = ``\uparrow", \forall s$
- $\gamma = 0.9$



$$Arr R(6,\uparrow) + \gamma [.2R(2,\uparrow) + .8R(3,\uparrow)]$$

$$= -10 + 0.9 * (0.2 * 0 + 0.8 * 1)$$

$$\Rightarrow -9.28$$

$$egin{aligned} \mathrm{R}(8,\uparrow) + \gamma \mathrm{R}(5,\uparrow) & \mathrm{R}(9,\uparrow) + \gamma \mathrm{R}(6,\uparrow) \ &= 0 + 0.9 * (-10) \Rightarrow -9 \end{aligned}$$



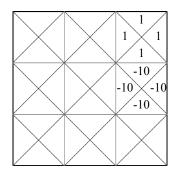
horizon h = 3

 $\gamma \mathrm{R}(2,\uparrow)$

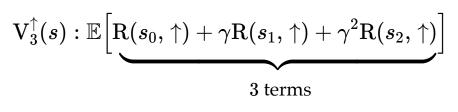
 $\gamma^2 \mathrm{R}(2,\uparrow)$

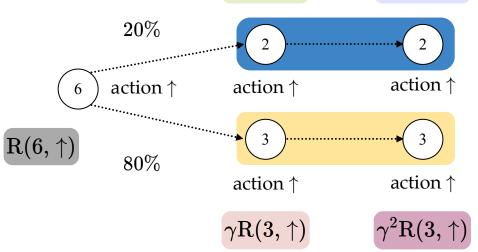
states and one special transition:

1	2 20% [▼] ···	3 • <u>*</u> 80%	
4	5	6	
7	8	9	



- $\pi(s) = ``\uparrow", \ \forall s$
- $\gamma = 0.9$





$$\begin{split} V_3^{\uparrow}(6) &= \mathbb{R}(6,\uparrow) + 20\% \quad \gamma \mathbb{R}(2,\uparrow) + 80\% \quad \gamma \mathbb{R}(3,\uparrow) + 20\% \quad \gamma^2 \mathbb{R}(2,\uparrow) + 80\% \quad \gamma^2 \mathbb{R}(3,\uparrow) \\ &= \mathbb{R}(6,\uparrow) + 20\% \quad [\quad \gamma \mathbb{R}(2,\uparrow) + \gamma^2 \mathbb{R}(2,\uparrow) \quad] \quad + 80\% \quad [\quad \gamma \mathbb{R}(3,\uparrow) \quad + \quad \gamma^2 \mathbb{R}(3,\uparrow) \quad] \\ &= \mathbb{R}(6,\uparrow) + 20\% \quad \gamma \quad [\mathbb{R}(2,\uparrow) + \gamma \mathbb{R}(2,\uparrow)] \quad + 80\% \quad \gamma \quad [\mathbb{R}(3,\uparrow) + \gamma \mathbb{R}(3,\uparrow)] \\ &= \mathbb{R}(6,\uparrow) + 20\% \quad \gamma \quad V_2^{\uparrow}(2) \quad + 80\% \quad \gamma \quad V_2^{\uparrow}(3) \end{split}$$

$$\mathrm{V}_3^{\uparrow}(6) = \mathrm{R}(6,\uparrow) \, + \, 20\% \,\, \gamma \,\, \mathrm{V}_2^{\uparrow}(2) \, + \, 80\% \,\, \gamma \,\, \mathrm{V}_2^{\uparrow}(3)$$

horizon-h value in state s: the expected sum of discounted rewards, starting in state s and following policy π for h steps.

(eq. 2)
$$V_h^{\pi}(s) = \mathrm{R}\left(s,\pi(s)\right) + \gamma \sum_{s'} \mathrm{T}\left(s,\pi(s),s'\right) V_{h-1}^{\pi}\left(s'\right)$$

the immediate reward for taking the policy-prescribed action $\pi(s)$ in state s.

(h-1) horizon future value at a next state s'

sum of future values weighted by the probability of reaching that next state s'

discounted by γ

Bellman Recursion (finite horizon h) $V_h^{\pi}(s) = R\left(s,\pi(s)\right) + \gamma \sum_{s'} T\left(s,\pi(s),s'\right) V_{h-1}^{\pi}\left(s'\right)$

states and one special transition:

1	2 20% [▼] ···.	3 ≜ 80%	
4	5	6	
7	8	9	

0 /	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	1
0 /	0 /	-10/
0×0	0×0	-10 <-10
0	0	-10
0 /	0 /	0 /
0×0	0×0	$ 0 \times 0 $
0	/ 0	0

- $\pi(s) = ``\uparrow", \ \forall s$
- $\gamma=0.9$

${ m V}_{1}^{\uparrow}$ ($(s) = \mathrm{I}$	$\mathrm{R}(s,\uparrow)$

0.00	0.00	1.00
0.00	0.00	-10.00
0.00	0.00	0.00

$${
m V}_2^{\uparrow}(s)$$

	2 ()	
0.00	0.00	1.90
0.00	0.00	-9.28
0.00	0.00	-9.00

$$egin{aligned} \mathbf{V}_2^\uparrow(9) &= \mathbf{R}(9,\uparrow) + \gamma [\mathbf{V}_1^\uparrow(6)] \ &= 0 + 0.9 imes [-10] \end{aligned}$$

$$= -9$$

Bellman Recursion (finite horizon h) $V_h^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^{\pi}(s')$

states and one special transition:

1	2	3	
4	20 % ···	80% 16	
7	8	9	

0 /	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	1
0 /	0 /	-10/
0×0	0×0	-10 <-10
0	0	-10
0 /	0 /	0 /
0×0	0×0	$\mid 0 \times 0 \mid$
0	0	0

- $\pi(s) = ``\uparrow", \forall s$
- $\gamma = 0.9$

${ m V}_{1}^{\uparrow}$ ((s) =	$=\mathrm{R}(s$	$,\uparrow)$

0.00	0.00	1.00
0.00	0.00	-10.00
0.00	0.00	0.00

${\rm V}_2^{\uparrow}(s)$			
0.00	0.00	1.90	
0.00	0.00	-9.28	
0.00	0.00	-9.00	

${ m V}_3^{\uparrow}(s)$			
0.00	0.00	2.71	
0.00	0.00	-8.63	
0.00	0.00	-8.35	

$$egin{aligned} \mathbf{V}_3^{\uparrow}(6) &= \mathbf{R}(6,\uparrow) + \gamma[.2 imes \mathbf{V}_2^{\uparrow}(2) + .8 imes \mathbf{V}_2^{\uparrow}(3)] \ &= -10 + .9[.2 imes \mathbf{0} + 0.8 imes \mathbf{1.9}] \ &= -8.632 \end{aligned}$$

Bellman Recursion (finite horizon h) $V_h^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^{\pi}(s')$

states and one special transition:

1	2 20% [▼] ···.	3 <u></u> *80%	
4	5	6	
7	8	9	

0 /	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	1
0 /	0 /	-10/
0×0	0×0	-10 \(-10
0	0	-10
0 /	0	0 /
0×0	0×0	$\mid 0 \times 0 \mid$
0	0	0

- $\pi(s) = ``\uparrow", \forall s$
- $\gamma = 0.9$

${ m V}_4^{\uparrow}(s)$			
0.00	0.00	3.44	
0.00	0.00	-8.05	
0.00	0.00	-7.77	

${ m V}_5^{\uparrow}(s)$		
0.00	0.00	4.10
0.00	0.00	-7.52
0.00	0.00	-7.24

	${ m V}_6^{\uparrow}(s)$		
0.00	0.00	4.69	
0.00	0.00	-7.05	• •
0.00	0.00	-6.77	

$$egin{aligned} \mathbf{V}_6^{\uparrow}(6) &= \mathbf{R}(6,\uparrow) + \gamma[.2 imes \mathbf{V}_5^{\uparrow}(2) + .8 imes \mathbf{V}_5^{\uparrow}(3)] \ &= -10 + .9[.2 imes \mathbf{0} + 0.8 imes 4.10] \ &= -7.048 \end{aligned}$$

Bellman Recursion (finite horizon h) $V_h^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^{\pi}(s')$

states and one special transition:

1	2 20% ···.	3 ♣ 80%
4	5	6
7	8	9

0 /	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	/1
0 /	0 /	-10/
0×0	0×0	-10 \(-10 \)
0	0	-10
0 /	0	0
0×0	0×0	$ 0 \times 0 $
0	0	0

- $\pi(s) = ``\uparrow", \forall s$
- $\gamma=0.9$

${ m V}_{59}^{\uparrow}(s)$		
0.00	0.00	9.98
0.00	0.00	-2.82
0.00	0.00	-2.54

${ m V}_{60}^{\uparrow}(s)$		
0.00	0.00	9.98
0.00	0.00	-2.81
0.00	0.00	-2.53

	$V_{61}^+(s)$		
0.00	0.00	9.98	
0.00	0.00	-2.81	• •
0.00	0.00	-2.53	

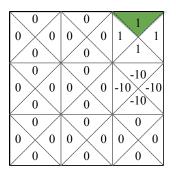
$$egin{aligned} \mathbf{V}_{61}^{\uparrow}(6) &= \mathrm{R}(6,\uparrow) + \gamma [.2 imes \mathbf{V}_{60}^{\uparrow}(2) + .8 imes \mathbf{V}_{60}^{\uparrow}(3)] \ &= -10 + .9 [.2 imes \mathbf{0} + 0.8 imes 9.98] \ &= -2.8144 \end{aligned}$$

Value functions converge as $h \to \infty$

states and one special transition:

1	2 20% ···.	3 ≜ 80%	
4	5	6	
7	8	9	

rewards



- $\pi(s) = ``\uparrow", \forall s$
- $\gamma = 0.9$

${ m V}_{60}^{\uparrow}(s)$		
0.00	0.00	9.98
0.00	0.00	-2.81
0.00	0.00	-2.53

${ m V}_{61}^{\uparrow}(s)$		
0.00	0.00	9.98
0.00	0.00	-2.81
0.00	0.00	-2.53

$\mathrm{V}^{\uparrow}_{\infty}(s)$		
0.00	0.00	10.00
0.00	0.00	-2.80
0.00	0.00	-2.52

- As we extend the horizon, value differences shrink
- because longer-term rewards are heavily discounted
- so, as $h \to \infty$, the value functions stop changing
- convergence can be seen, e.g., via $V_\infty^\uparrow(3)=1+.9+.9^2+.9^3+\cdots=10$

Typically, $\gamma < 1$ to ensure V_{∞} is finite.

As horizon $h \to \infty$, the Bellman recursion becomes the Bellman equation

states and one special transition:

1	2 20% [▼] ···.	3 ♣80%
4	5	6
7	8	9

rewards

0 /	0 /	1
0×0	0×0	1 1
0	0	1
0 /	0 /	-10/
0×0	0×0	-10 \(-10 \)
0	0	-10
0 /	0 /	0
0×0	0×0	0×0
/ 0 \		0

- $\pi(s) = ``\uparrow", \forall s$
- $\gamma=0.9$

Recursion (finite h) 2 $V_h^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^{\pi}(s')$

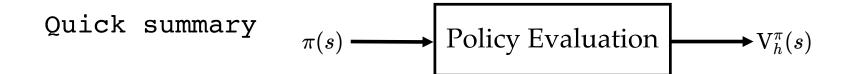
Equation
$$(h o \infty)$$
 3 $V^\pi_\infty(s) = \mathrm{R}(s,\pi(s)) + \gamma \sum_{s'} \mathrm{T}\left(s,\pi(s),s'\right) V^\pi_\infty(s')$

A system of |S| self-consistent linear equations, one for each state

$\mathrm{V}^{\uparrow}_{\infty}(s)$				
0.00	0.00	10.00		
0.00	0.00	-2.80		
0.00	0.00	-2.52		

$$egin{aligned} \mathbf{V}_{\infty}^{\uparrow}(3) &= \mathbf{R}(3,\uparrow) + \gamma [\mathbf{V}_{\infty}^{\uparrow}(3)] \ &= 1 + .9 imes 10 \Rightarrow 10 \end{aligned}$$

$$egin{aligned} \mathbf{V}_{\infty}^{\uparrow}(6) &= \mathrm{R}(6,\uparrow) + \gamma[.2 imes \mathbf{V}_{\infty}^{\uparrow}(2) + .8 imes \mathbf{V}_{\infty}^{\uparrow}(3)] \ &= -10 + .9[.2 imes \mathbf{0} + 0.8 imes \mathbf{10}] \Rightarrow -2.8 \end{aligned}$$



Use the definition and sum up expected rewards:

$$\mathbf{V}_h^{\pi}(s) := \mathbb{E}\left[\sum_{t=0}^{h-1} \gamma^t \mathrm{R}\left(s_t, \pi\left(s_t
ight)
ight) \mid s_0 = s, \pi
ight]$$

Or, leverage the recursive structure:

2 finite-horizon Bellman recursions

$$\mathrm{V}_h^{\pi}(s) = \mathrm{R}(s,\pi(s)) + \gamma \sum_{s'} \mathrm{T}\left(s,\pi(s),s'
ight) \mathrm{V}_{h-1}^{\pi}\left(s'
ight)$$

infinite-horizon Bellman equations

$$\mathrm{V}_{\infty}^{\pi}(s) = \mathrm{R}(s,\pi(s)) + \gamma \sum_{s'} \mathrm{T}\left(s,\pi(s),s'
ight) \mathrm{V}_{\infty}^{\pi}\left(s'
ight)$$

Outline

- Markov Decision Processes Definition
- Policy Evaluation
 - State value functions: V^{π}
 - Bellman recursions and Bellman equations
- Policy Optimization
 - Optimal policies π^*
 - Optimal action value functions: Q*
 - Value iteration

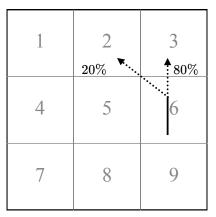
Optimal policy π^*

- Intuitively, an optimal policy π^* is a policy that yields the highest possible value $V_h^*(s)$ from every state
- An MDP has a unique optimal value $\mathrm{V}_h^*(s)$
- Optimal policy π^* might not be unique

e.g. in the "Luigi game", all rewards are 1,



States and one special transition:



Rewards:

:	1	1	1
	1 💢 1	1 × 1	1 💢 1
	1	1	_ 1
	1	1	1
	1 1	1 × 1	1 💢 1
	1	1	1
	1	1	1
	1 1	1 🔨 1	1 × 1
	1	1	1

 $\gamma = 0.9$

then any policy is an optimal policy

Optimal policy π^*

- Formally: an optimal policy π^* is such that: $V_h^{\pi^*}(s) = \max_{\pi} V_h^{\pi}(s) = V_h^*(s), \forall s \in \mathcal{S}$
- How to search for an optimal policy π^* ?
- Even if we tediously enumerate over all π , do policy evaluation, compare values to get $V_h^*(s)$...it's not yet clear how to choose actions.

 $V^*(s)$ is defined over states, not actions.

It tells us where we'd like to *be* — not what we should do to *get* there.

Optimality recursion

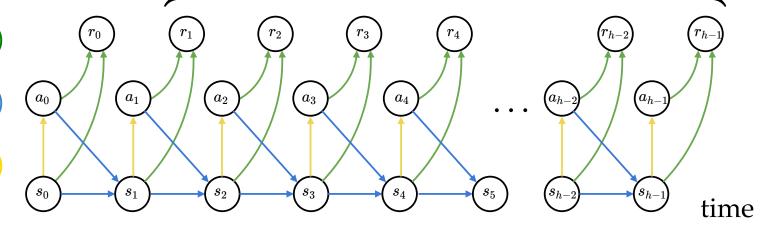
if we've acted optimally for h steps: $V_h^*(s)$

we must have acted optimally from the first step onward $V_{h-1}^*(s')$

Reward R(s, a)

Transition T (s, a, s')

Policy $\pi(s)$



Bellman recursion under an optimal policy

$$\mathrm{V}_h^*(s) = \mathrm{max}_a \left[\mathrm{R}(s,a) + \gamma \sum_{s'} \mathrm{T}(s,a,s') \mathrm{V}_{h-1}^*(s')
ight]$$

Define the optimal state-action value functions $Q_h^*(s, a)$:

the expected sum of discounted rewards, obtained by

- starting in state *s*
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

$$egin{align} egin{align} egin{align} \mathbf{V}_h^*(s) &= \max_a \left[\mathbf{R}(s,a) + \gamma \sum_{s'} \mathbf{T}(s,a,s') \mathbf{V}_{h-1}^*(s')
ight] &= \max_a \left[\mathbf{Q}_h^*(s,a)
ight] \ \mathbf{Q}_h^*(s,a) \end{aligned}$$

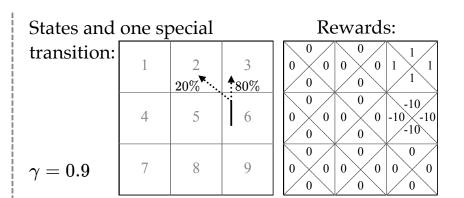
Q* satisfies the Bellman recursion:

$$\mathrm{Q}_{h}^{st}(s,a) = \mathrm{R}(s,a) + \gamma \sum_{s'} \mathrm{T}\left(s,a,s'
ight) \mathrm{max}_{a'} \, \mathrm{Q}_{h-1}^{st}\left(s',a'
ight)$$

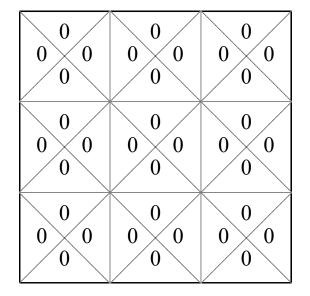


 $\mathbf{Q}_h^*(s,a)$: the value for

- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps



$$\mathrm{Q}_0^*(s,a)$$



$$\mathrm{Q}_1^*(s,a) = \mathrm{R}(s,a)$$

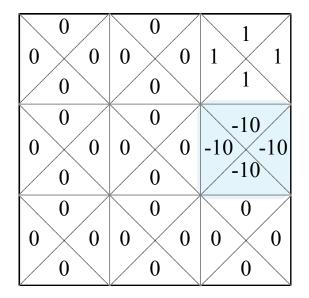
	0			0		$\setminus 1$
0	X	0	0	\times	0	$ 1 \times 1 $
	0			0		1
	0			0		-10/
0	\times	0	0	\times	0	-10 < -10
	0			0		-10
	0			0	//	0 /
0	X	0	0	\times	0	$\mid 0 \searrow 0 \mid$
	0			0		0



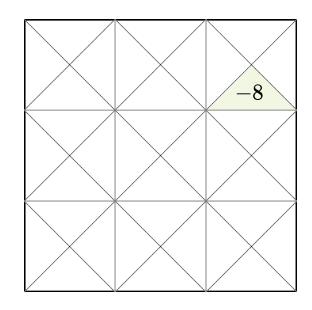
$Q_h^*(s, a)$: the value for

- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

$$\mathrm{Q}_1^*(s,a)$$



$$\mathrm{Q}_2^*(s,a)$$



Consider $Q_2^*(3,\downarrow)$

- receive $R(3,\downarrow)$
- next state s' = 6, act **optimally** for the remaining one timestep
 - receive $\max_{a'} \mathbf{Q}_1^* (6, a')$

$$egin{aligned} \mathbf{Q}_2^*(\mathbf{3},\downarrow) &= \mathbf{R}(\mathbf{3},\downarrow) \ + \gamma \max_{a'} \mathbf{Q}_1^*\left(\mathbf{6},a'
ight) \ &= 1+.9 imes 10 \ &= -8 \end{aligned}$$

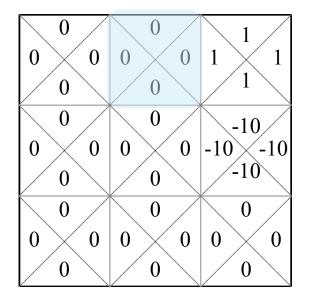


$\mathrm{Q}_h^*(s,a)$: the value for

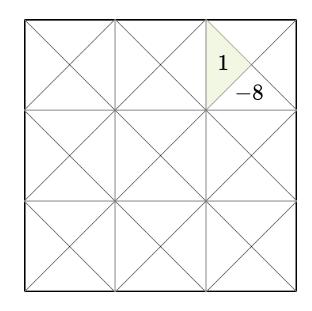
- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

States and	one sp	ecial		Rewards:
transition:	1	2	3	
	4	20%*···	1 80% 6	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$
$\gamma=0.9$	7	8	9	0 0 0 0 0 0

$$\mathrm{Q}_1^*(s,a)$$



$$\mathrm{Q}_2^*(s,a)$$



Let's consider $Q_2^*(3, \leftarrow)$

- receive $R(3, \leftarrow)$
- next state s' = 2, act **optimally** for the remaining one timestep
 - receive $\max_{a'} \mathrm{Q}_1^* \left(2, a' \right)$

$$egin{aligned} \mathbf{Q}_2^*(\mathbf{3},\leftarrow) &= \mathbf{R}(\mathbf{3},\leftarrow) \,+ \gamma \max_{a'} \mathbf{Q}_1^*\left(\mathbf{2},a'
ight) \ &= \mathbf{1} + .9 imes 0 \ &= \mathbf{1} \end{aligned}$$



$Q_h^*(s, a)$: the value for

- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

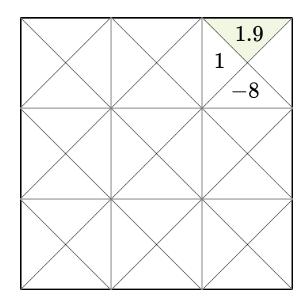
States and one special transition:

1.0 01	ric special					
1	2 20% ···.	3 •80%				
4	5	6				
7	8	9				

 $\mathrm{Q}_1^*(s,a)$

0	0 /	1
0×0	0×0	$ 1 \times 1 $
0	0	1
0	0 /	-10/
0×0	0×0	-10 \(\sigma -10 \)
0	0	-10
0 /	0 /	\setminus 0 $/$
0×0	0×0	$\mid 0 \times 0 \mid$
0	0	0

 $\mathrm{Q}_2^*(s,a)$



Let's consider $Q_2^*(3,\uparrow)$

 $\gamma = 0.9$

- receive $R(3,\uparrow)$
- next state s' = 3, act **optimally** for the remaining one timestep
 - receive $\max_{a'} \mathbf{Q}_1^* (3, a')$

$$egin{aligned} \mathbf{Q}_2^*(\mathbf{3},\uparrow) &= \mathbf{R}(\mathbf{3},\uparrow) \ + \gamma \max_{a'} \mathbf{Q}_1^*\left(\mathbf{3},a'
ight) \ &= \mathbf{1} + .9 imes \mathbf{1} \ &= \mathbf{1}.9 \end{aligned}$$



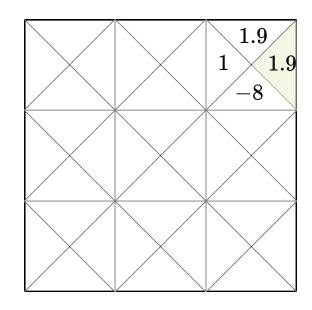
$Q_h^*(s, a)$: the value for

- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

$$\mathrm{Q}_1^*(s,a)$$

\setminus 0 /	0 /	1
0×0	0×0	1×1
0	0	/1
\setminus 0 /	0 /	-10/
0×0	0×0	-10 < -10
0	0	-10
\setminus 0 /	0 /	\setminus 0 $/$
0×0	0×0	$\mid 0 \times 0 \mid$
0	0	0

$$\mathrm{Q}_2^*(s,a)$$



Let's consider $Q_2^*(3, \rightarrow)$

- receive $R(3, \rightarrow)$
- next state s' = 3, act **optimally** for the remaining one timestep
 - receive $\max_{a'} \mathbf{Q}_1^* (3, a')$

$$egin{aligned} \mathbf{Q}_2^*(3, o) &= \mathbf{R}(3, o) \, + \gamma \max_{a'} \mathbf{Q}_1^* \, (3, a') \ &= 1 + .9 imes 1 \ &= 1.9 \end{aligned}$$



$\mathrm{Q}_h^*(s,a)$: the value for

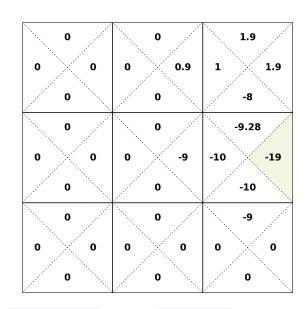
- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

States and	one sp	Rewards:		
transition:	1	2 20% ▼ ···.	3 • 80%	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$
	4	5	6	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
$\gamma=0.9$	7	8	9	

$$\mathrm{Q}_1^*(s,a)$$

0	0	1
0 0	0 0	1 1
0	0	1
0	0	-10
0 0	0 0	-10
0	0	-10
0	0	0
0 0	0 0	0 0
0	0	0

$$\mathrm{Q}_2^*(s,a)$$



Let's consider $Q_2^*(6, \rightarrow)$

- receive $R(6, \rightarrow)$
- act optimally at the next state s'=6 receive $\max_{a'} \mathrm{Q}_1^* \left(6, a' \right)$

$$egin{aligned} \mathbf{Q}_2^*(6,
ightarrow) &= \mathrm{R}(6,
ightarrow) + \gamma [\mathrm{max}_{a'} \, \mathrm{Q}_1^* \, (6,a')] \ &= -10 + .9 imes -10 \Rightarrow -19 \end{aligned}$$



$Q_h^*(s, a)$: the value for

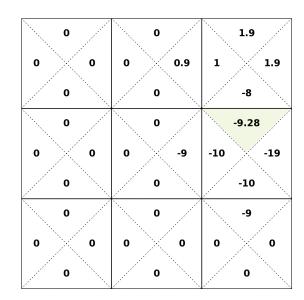
- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

States and	one sp	pecial		Rewards:
transition:	1	2 20% ▼ ···.	3 • 80%	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$
	4	5	6	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
$\gamma=0.9$	7	8	9	$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 &$

$$\mathrm{Q}_1^*(s,a)$$

0	0	1
0 0	0 0	1 1
0	0	1
0	0	-10
0 0	0 0	-10 -10
0	0	-10
0	0	0
0 0	0 0	0 0
0	0	0

$$\mathrm{Q}_2^*(s,a)$$



Let's consider $Q_2^*(6,\uparrow)$

- receive $R(6,\uparrow)$
- act **optimally** at the next state s'
 - 20% chance, s' = 2, act optimally, get $\max_{a'} \mathbf{Q}_1^* \left(2, a' \right)$
 - 80% chance, s' = 3, act optimally, get $\max_{a'} \mathbf{Q}_1^* (3, a')$

$$egin{aligned} \mathbf{Q}_2^*(6,\uparrow) &= \mathbf{R}(6,\uparrow) \ + \gamma [.2 \max_{a'} \mathbf{Q}_1^* \, (2,a') + .8 \max_{a'} \mathbf{Q}_1^* \, (3,a')] \ &= -10 + .9 [.2 imes 0 + .8 imes 1] \Rightarrow -9.28 \end{aligned}$$

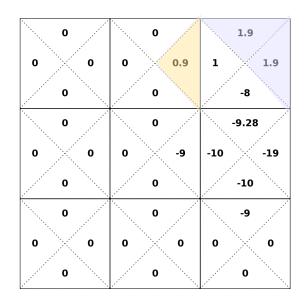


$Q_h^*(s, a)$: the value for

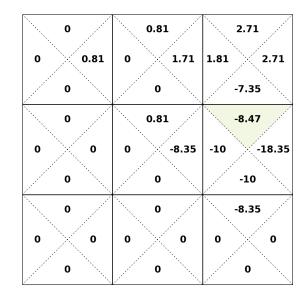
- starting in state *s*,
- take action *a*, for one step
- act **optimally** thereafter for the remaining (h-1) steps

States and one special			Rewards:	
transition:	1	2 20%*··	3	
	4	20%*···	1 80% 6	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
$\gamma=0.9$	7	8	9	

$$\mathrm{Q}_2^*(s,a)$$



$$\mathrm{Q}_3^*(s,a)$$



Let's consider $Q_3^*(6,\uparrow)$

- receive $R(6,\uparrow)$
- act **optimally** at the next state s'
 - 20% chance, s' = 2, act optimally, get $\max_{a'} \mathbf{Q}_2^* \left(2, a' \right)$
 - 80% chance, s' = 3, act optimally, get $\max_{a'} \mathbf{Q}_2^* (3, a')$

$$egin{aligned} \mathbf{Q}_{3}^{*}(6,\uparrow) &= \mathrm{R}(6,\uparrow) \ + \gamma[.2\max_{a'}\mathbf{Q}_{2}^{*}\left(2,a'
ight) + .8\max_{a'}\mathbf{Q}_{2}^{*}\left(3,a'
ight)] \ &= -10 + .9[.2 imes rac{0.9}{0.9} + .8 imes 1.9] \Rightarrow -8.47 \end{aligned}$$

Value iteration: what we just did, iteratively invoke 5

$$\mathrm{Q}_{h}^{st}(s,a) = \mathrm{R}(s,a) + \gamma \sum_{s'} \mathrm{T}\left(s,a,s'
ight) \max_{a'} \mathrm{Q}_{h-1}^{st}\left(s',a'
ight)$$

Value Iteration

1. for
$$s \in \mathcal{S}, a \in \mathcal{A}$$
:

2.
$$Q_{old}(s, a) = 0$$

3. **while** True:

if run this block *h* times

$$4. \quad ext{ for } s \in \mathcal{S}, a \in \mathcal{A}:$$

frun this block
$$h$$
 times and break, then the returns are exactly Q_h^* $\begin{cases} 4. & \text{for } s \in \mathcal{S}, a \in \mathcal{A}: \\ Q_{\text{new}}\left(s,a\right) \leftarrow \mathbb{R}(s,a) + \gamma \sum_{s'} \operatorname{T}\left(s,a,s'\right) \max_{a'} Q_{\text{old}}\left(s',a'\right) \end{cases}$ $\begin{cases} 6. & \text{if } \max_{s,a} |Q_{\text{old}}\left(s,a\right) - Q_{\text{new}}\left(s,a\right)| < \epsilon: \\ 7. & \text{return } Q_{\text{new}} \end{cases}$ $\begin{cases} 8. & Q_{\text{old}} \leftarrow Q_{\text{new}} \end{cases}$

$$\mathbf{if} \ \mathrm{max}_{s,a} \left| Q_{\mathrm{old}} \left(s,a \right) - Q_{\mathrm{new}} \left(s,a \right) \right| < \epsilon$$

$$\mathrm{Q}^*_\infty(s,a)$$

$$\mathrm{Q}_{\mathrm{old}} \, \leftarrow \mathrm{Q}_{\mathrm{new}}$$

Optimal policy easily extracted: 6 $\pi_h^*(s) = rg \max_a \mathrm{Q}_h^*(s,a)$

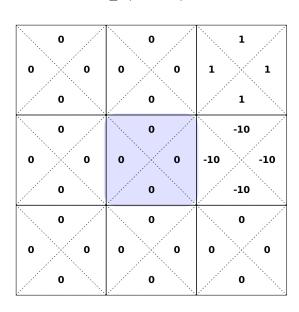
e.g. the best actions to take in state 5

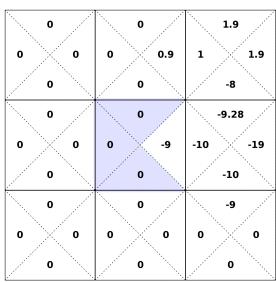
$$\mathrm{Q}_1^*(s,a)$$

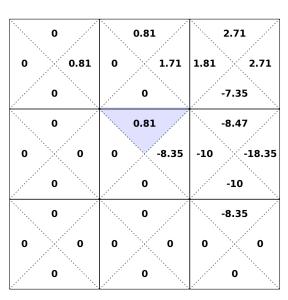
$$\mathrm{Q}_2^*(s,a)$$

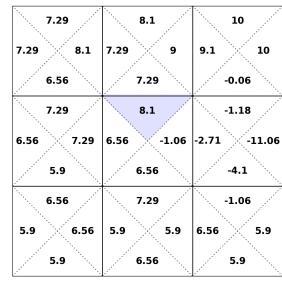
$$Q_3^*(s,a)$$

$$\mathrm{Q}_3^*(s,a) \qquad \qquad \ldots \qquad \mathrm{Q}_\infty^*(s,a)$$









- For finite h, optimal policy π_h^* depends on how many time steps are left
- When $h \to \infty$, time no longer matters, i.e., there exists a stationary π^*

Summary

- Markov decision processes (MDP) are a nice mathematical framework for making sequential decisions. It's the foundation to reinforcement learning.
- An MDP is defined by a five-tuple, and the goal is to find an optimal policy that leads to high expected cumulative discounted rewards.
- To evaluate how good a *given* policy π is, we can calculate $V^{\pi}(s)$ via
 - the summation-over-rewards definition
 - Bellman recursion for finite horizon and Bellman equation for infinite horizon
- To *find* an optimal policy, we can recursively find $Q^*(s, a)$ via the value iteration algorithm, and then act greedily w.r.t. the Q^* values.

We'd love to hear

your thoughts.
Thanks!