

<https://introml.mit.edu/>

6.390 Intro to Machine Learning

Lecture 10: Markov Decision Processes

Shen Shen

April 18, 2025

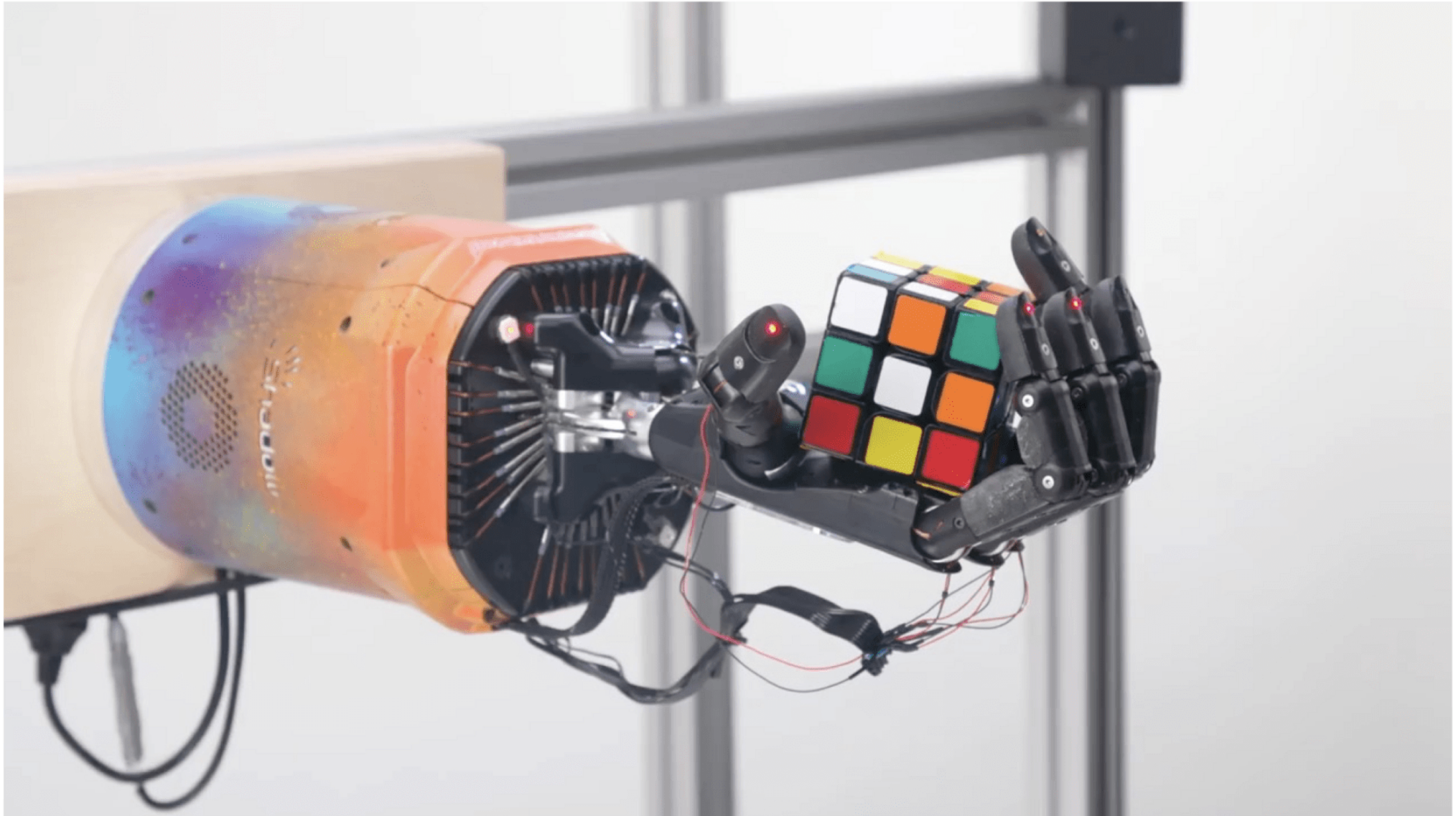
11am, Room 10-250

Learning to Walk

Massachusetts Institute of
Technology, 2004



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



[Hu

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

15]

Discovering faster matrix multiplication algorithms with reinforcement learning

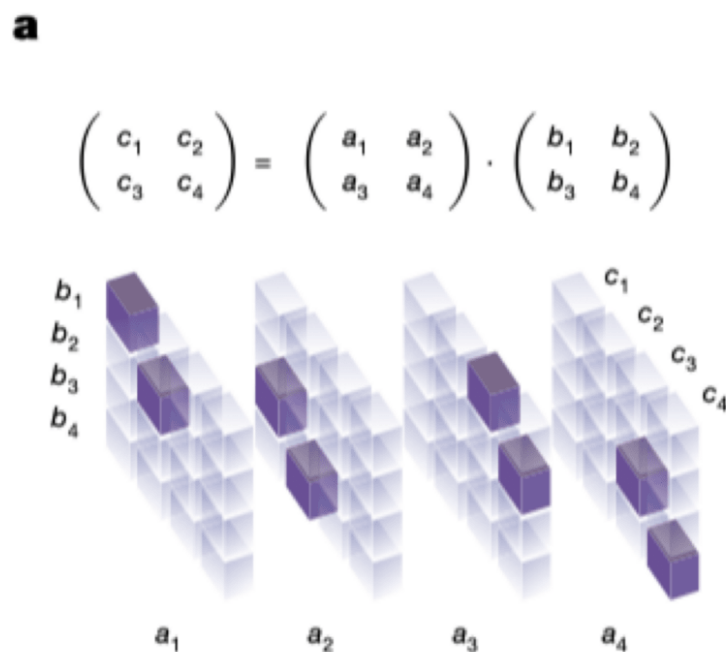
<https://doi.org/10.1038/s41586-022-05172-4>

Received: 2 October 2021

Accepted: 2 August 2022

Published online: 5 October 2022

Alhussein Fawzi^{1,2✉}, Matej Balog^{1,2}, Aja Huang^{1,2}, Thomas Hubert^{1,2}, Bernardino Romera-Paredes^{1,2}, Mohammadamin Barekatin¹, Alexandre Francisco J. R. Ruiz¹, Julian Schrittwieser¹, Grzegorz Swirszcz¹, David Si & Pushmeet Kohli¹



b

$$\begin{aligned} m_1 &= (a_1 + a_4)(b_1 + b_4) \\ m_2 &= (a_3 + a_4)b_1 \\ m_3 &= a_1(b_2 - b_4) \\ m_4 &= a_4(b_3 - b_1) \\ m_5 &= (a_1 + a_2)b_4 \\ m_6 &= (a_3 - a_1)(b_1 + b_2) \\ m_7 &= (a_2 - a_4)(b_3 + b_4) \\ c_1 &= m_1 + m_4 - m_5 + m_7 \\ c_2 &= m_3 + m_5 \\ c_3 &= m_2 + m_4 \\ c_4 &= m_1 - m_2 + m_3 + m_6 \end{aligned}$$

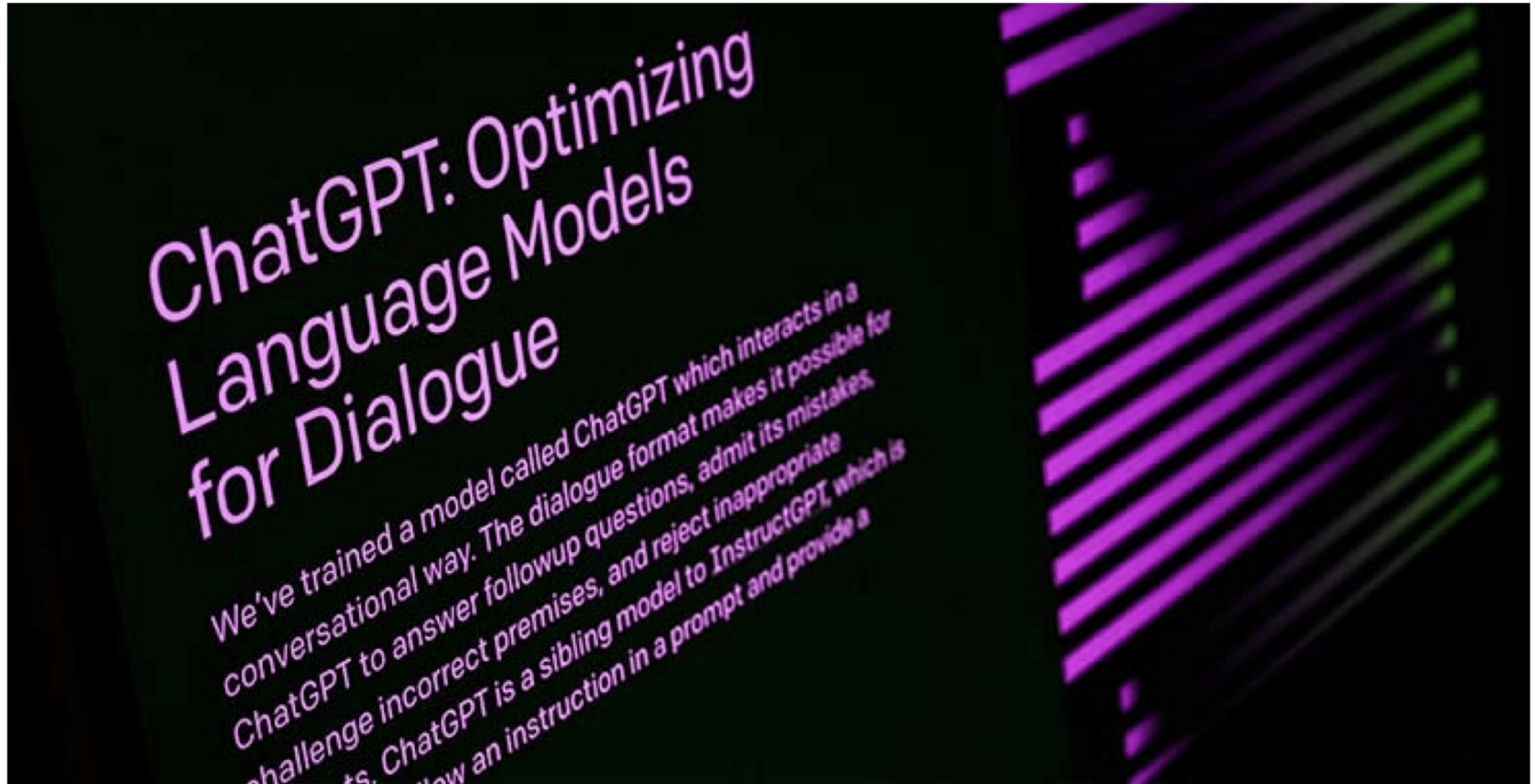
c

$$\begin{aligned} \mathbf{U} &= \begin{pmatrix} \\ \\ \\ \end{pmatrix} \\ \mathbf{V} &= \begin{pmatrix} \\ \\ \\ \end{pmatrix} \\ \mathbf{W} &= \begin{pmatrix} \\ \\ \\ \end{pmatrix} \end{aligned}$$

Size (n, m, p)	Best method known	Best rank known	AlphaTensor rank Modular	AlphaTensor rank Standard
(2, 2, 2)	(Strassen, 1969) ²	7	7	7
(3, 3, 3)	(Laderman, 1976) ¹⁵	23	23	23
(4, 4, 4)	(Strassen, 1969) ² (2, 2, 2) \otimes (2, 2, 2)	49	47	49
(5, 5, 5)	(3, 5, 5) + (2, 5, 5)	98	96	98
(2, 2, 3)	(2, 2, 2) + (2, 2, 1)	11	11	11
(2, 2, 4)	(2, 2, 2) + (2, 2, 2)	14	14	14
(2, 2, 5)	(2, 2, 2) + (2, 2, 3)	18	18	18
(2, 3, 3)	(Hopcroft and Kerr, 1971) ¹⁶	15	15	15
(2, 3, 4)	(Hopcroft and Kerr, 1971) ¹⁶	20	20	20
(2, 3, 5)	(Hopcroft and Kerr, 1971) ¹⁶	25	25	25
(2, 4, 4)	(Hopcroft and Kerr, 1971) ¹⁶	26	26	26
(2, 4, 5)	(Hopcroft and Kerr, 1971) ¹⁶	33	33	33
(2, 5, 5)	(Hopcroft and Kerr, 1971) ¹⁶	40	40	40
(3, 3, 4)	(Smirnov, 2013) ¹⁸	29	29	29
(3, 3, 5)	(Smirnov, 2013) ¹⁸	36	36	36
(3, 4, 4)	(Smirnov, 2013) ¹⁸	38	38	38
(3, 4, 5)	(Smirnov, 2013) ¹⁸	48	47	47
(3, 5, 5)	(Sedoglavic and Smirnov, 2021) ¹⁹	58	58	58
(4, 4, 5)	(4, 4, 2) + (4, 4, 3)	64	63	63
(4, 5, 5)	(2, 5, 5) \otimes (2, 1, 1)	80	76	76

X

Reinforcement Learning with Human Feedback



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

Outline

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

Outline

- Markov Decision Processes Definition, terminologies, and policy
- 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 accessible and elegant theory, math, algorithms, and applications. But also, considerable variation in notations.
- We will use the most RL-flavored notations.



Running example: Mario in a grid-world

1	2	3
4	5	6
7	8	9

Diagram illustrating a 3x3 grid world. The grid is numbered 1 to 9. From state 6 (row 2, column 3), an action 'up' leads to state 5 (row 2, column 2) with a 20% probability (indicated by a dotted arrow) and to state 3 (row 1, column 3) with an 80% probability (indicated by a solid arrow).

- 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 “ \uparrow ” gets to state (4)
 - If an action would take Mario out of the grid world, stay put;
 - E.g., in state (9), “ \rightarrow ” gets back to state (9)
 - In state (6), action “ \uparrow ” leads to two possibilities:
 - 20% chance to (2)
 - 80% chance to (3).



Mario in a grid-world, cont'd

- (state, action) pairs give **rewards**:



- in state 3, any action gives reward 1

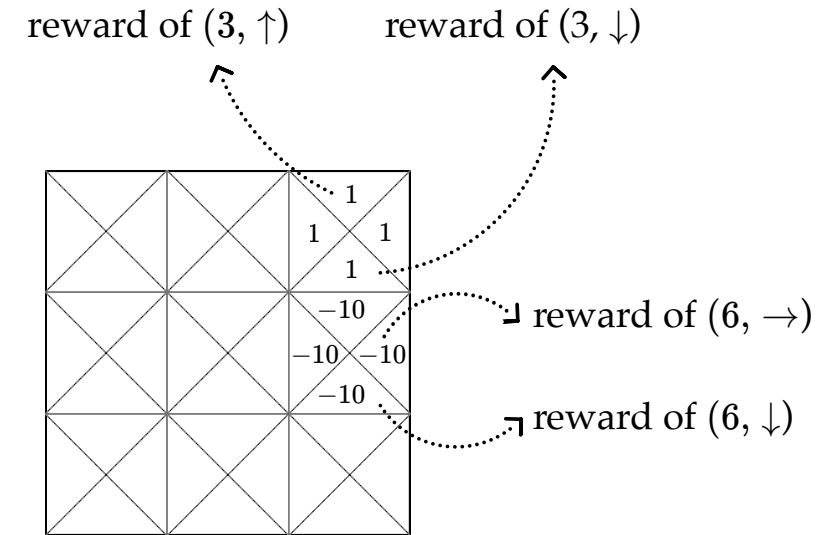


- in state 6, any action gives reward -10



- any other (state, action) pair gives reward 0

- **discount factor**: a scalar that reduces the "worth" of rewards, depending on the timing Mario gets the rewards.
 - e.g., say this factor is 0.9. then, for $(3, \leftarrow)$ pair, Mario gets a reward of 1 at the start of the game; at the 2nd time step, a discounted reward of 0.9; at the 3rd time step, it is further discounted to $(0.9)^2$, and so on.



Markov Decision Processes - Definition and terminologies

- \mathcal{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.

Markov Decision Processes - Definition and terminologies

- \mathcal{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
4	5	6
7	8	9

$$T(7, \uparrow, 4) = 1$$

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

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

$$T(6, \uparrow, 2) = 0.2$$

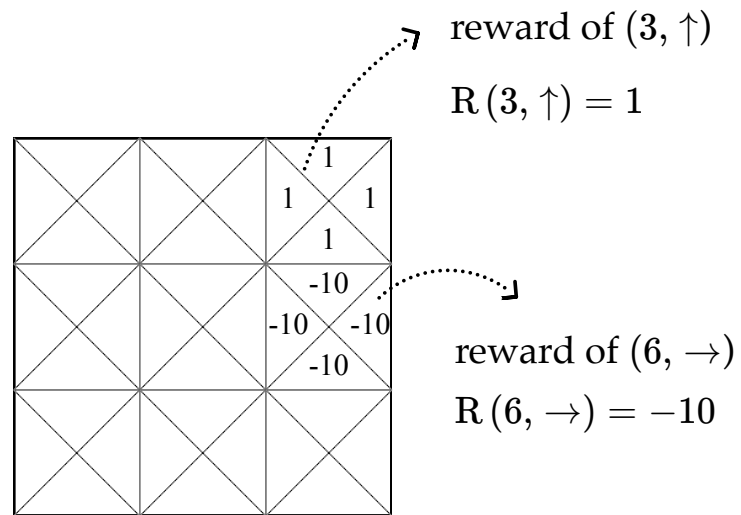
In 6.390,

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



Markov Decision Processes - Definition and terminologies

- \mathcal{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
- $R(s, a)$ is deterministic and bounded.

Markov Decision Processes - Definition and terminologies

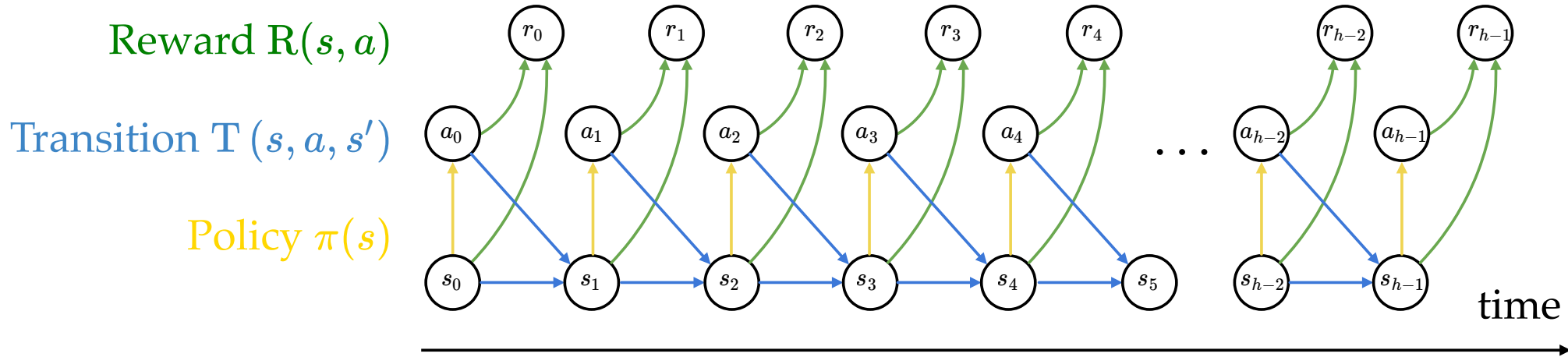
- \mathcal{S} : state space, contains all possible states s .
- \mathcal{A} : action space, contains all possible actions a .
- $T(s, a, s')$: the probability of transition from state s to s' when action a is taken.
- $R(s, a)$: reward, takes in a (state, action) pair and returns a reward.
- $\gamma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



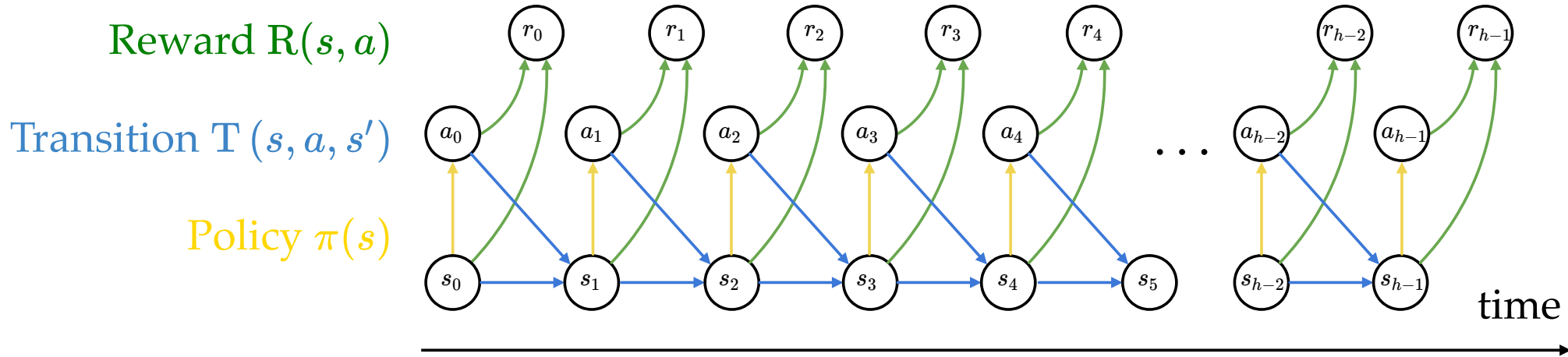
a trajectory (aka, an experience, or a rollout), of horizon h

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{h-1}, a_{h-1}, r_{h-1})$$

initial state

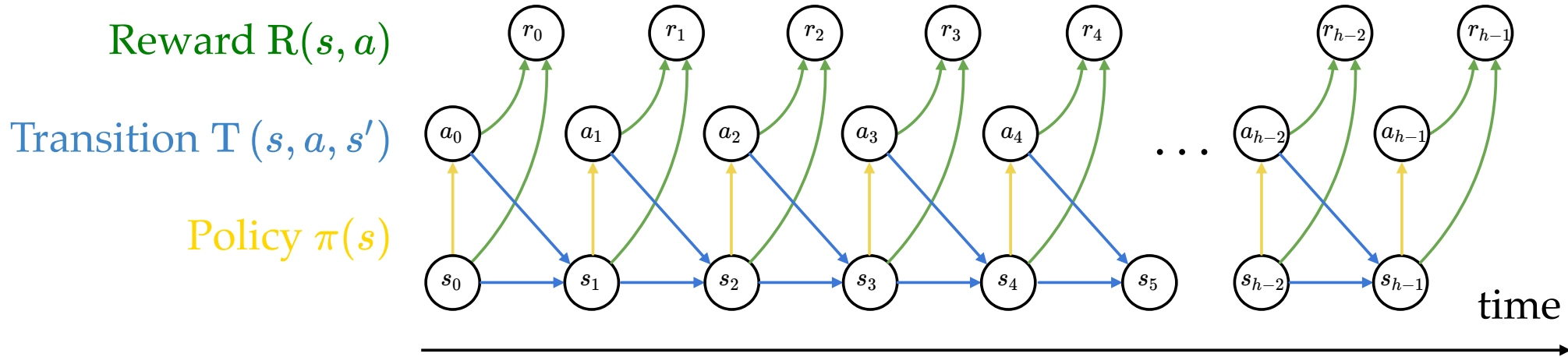
all depends on π

- $a_t = \pi(s_t)$
- $r_t = \mathbf{R}(s_t, a_t)$
- $\Pr(s_t = s' \mid s_{t-1} = s, a_{t-1} = a) = \mathbf{T}(s, a, s')$



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

One idea: $R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \gamma^3 R(s_3, a_3) + \dots + \gamma^{h-1} R(s_{h-1}, a_{h-1})$

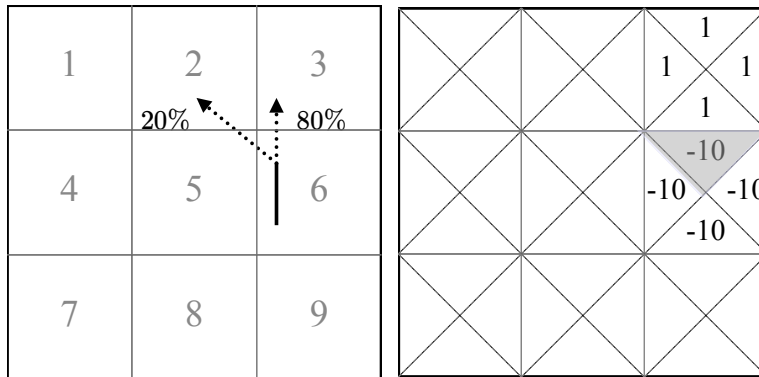


Starting in a given s_0 , how "good" is it to follow a policy π for h time steps?

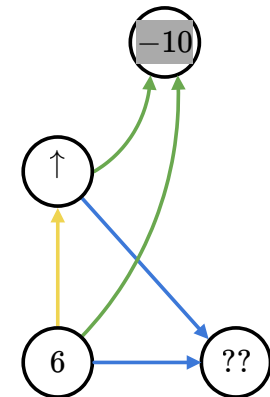
One idea: $R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \gamma^3 R(s_3, a_3) + \dots + \gamma^{h-1} R(s_{h-1}, a_{h-1})$

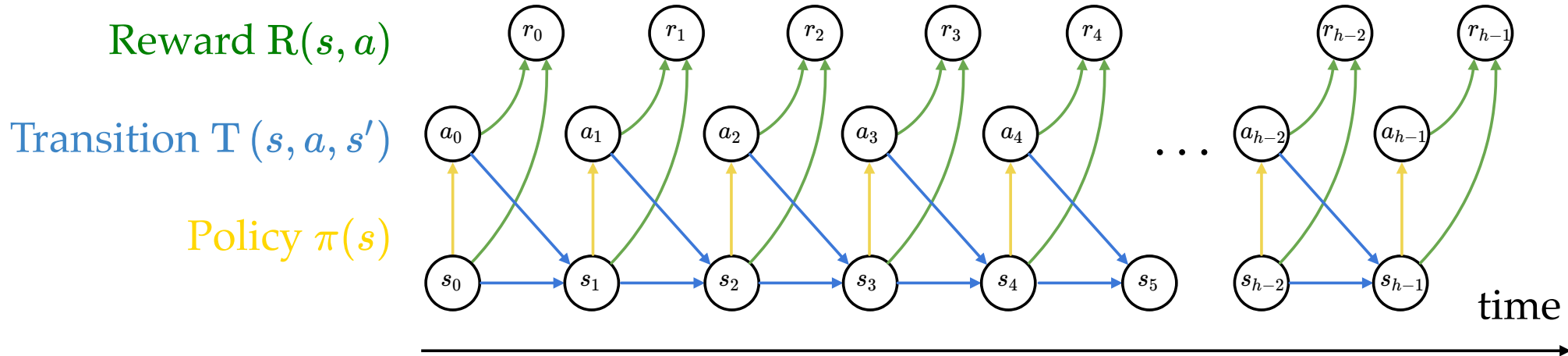
But in

Mario game:



if start at $s_0 = 6$ and policy $\pi(s) = \uparrow, \forall s$, i.e., always up





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

$$\mathbb{E} \left[\overbrace{R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \gamma^3 R(s_3, a_3) + \dots + \gamma^{h-1} R(s_{h-1}, a_{h-1})}^{h \text{ terms}} \right]$$

in 390, this expectation is only w.r.t. the transition probabilities $T(s, a, s')$

Outline

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

$$\mathbb{E}\left[\mathbf{R}(s_0, a_0) + \gamma\mathbf{R}(s_1, a_1) + \gamma^2\mathbf{R}(s_2, a_2) + \gamma^3\mathbf{R}(s_3, a_3) + \dots + \gamma^{h-1}\mathbf{R}(s_{h-1}, a_{h-1})\right]$$

Definition: For a *given* policy $\pi(s)$, the state **value functions**

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

- value functions $V_h^\pi(s)$: the expected sum of discounted rewards, starting in state s , and follow policy π for h steps.
- horizon-0 values defined as 0.
- value is long-term, reward is short-term (one-time).



evaluate the " $\pi(s) = \uparrow$, for all s , i.e. the always \uparrow " policy

states and
one special transition:

1	2	3
4	5	6
7	8	9

rewards

		1
	1	1
	1	
	-10	
	-10	-10
		-10

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

expanded form

$$\mathbb{E} \left[\underbrace{R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots}_{h \text{ terms}} \right]$$

h terms

horizon $h = 0$: no step left

$$V_0^\uparrow(s) = 0$$

0	0	0
0	0	0
0	0	0

horizon $h = 1$: receive the rewards

$$V_1^\uparrow(s) = R(s, \uparrow)$$

0	0	1
0	0	-10
0	0	0



horizon $h = 2$:

states and
one special transition:

1	2	3
4	5	6
7	8	9

Transitions from state 2: 20% to state 3, 80% to state 5.

rewards

		1
		1
		1
		-10
		-10
		-10

- $\pi(s) = \text{“}\uparrow\text{”}, \forall s$
- $\gamma = 0.9$

$$V_2^\uparrow(s) = \underbrace{\mathbb{E} \left[R(s_0, a_0) + \gamma R(s_1, a_1) \right]}_{\text{2 terms inside}}$$

$$R(1, \uparrow) + \gamma R(1, \uparrow)$$

$$R(2, \uparrow) + \gamma R(2, \uparrow)$$

$$R(3, \uparrow) + \gamma R(3, \uparrow)$$

$$R(4, \uparrow) + \gamma R(1, \uparrow)$$

0	0	1.9
0	0	

$$= 1 + 0.9 * (1) = 1.9$$

$$R(5, \uparrow) + \gamma R(2, \uparrow)$$



horizon $h = 2$:

states and
one special transition:

1	2	3
4	5	6
7	8	9

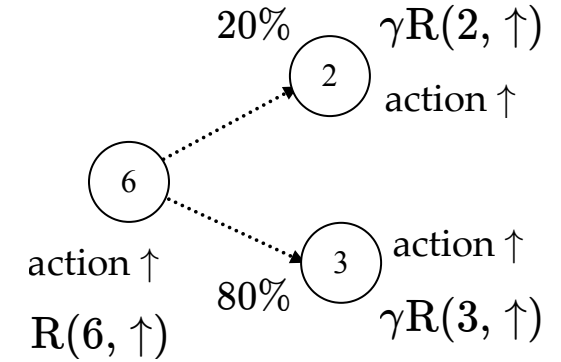
Transitions: 2 to 3 (80%), 2 to 4 (20%)

rewards

		1
		1
		1
		-10
		-10
		-10

- $\pi(s) = \text{“}\uparrow\text{”}, \forall s$
- $\gamma = 0.9$

$$V_2^\uparrow(s) = \underbrace{\mathbb{E} \left[R(s_0, a_0) + \gamma R(s_1, a_1) \right]}_{\text{2 terms inside}}$$



0	0	1.9
0	0	-9.28
0	0	-9

$$\begin{aligned} & \rightarrow R(6, \uparrow) + \gamma[.2R(2, \uparrow) + .8R(3, \uparrow)] \\ & = -10 + 0.9 * (0.2 * 0 + 0.8 * 1) \\ & = -9.28 \end{aligned}$$

$$R(7, \uparrow) + \gamma R(4, \uparrow)$$

$$R(8, \uparrow) + \gamma R(5, \uparrow)$$

$$R(9, \uparrow) + \gamma R(6, \uparrow) = 0 + 0.9 * (-10)$$



horizon $h = 3$: $\mathbb{E} [R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2)]$

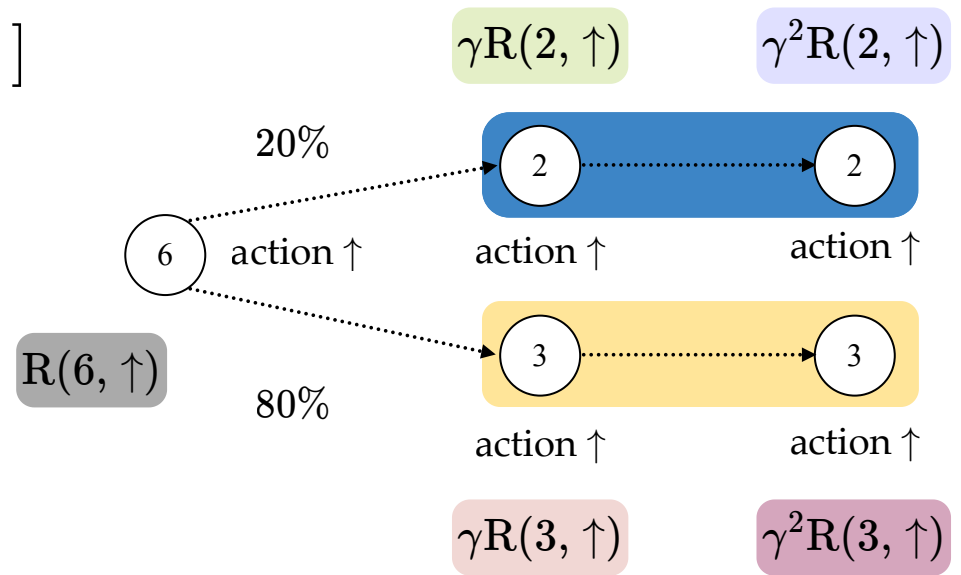
states and
one special transition:

1	2	3
4	5	6
7	8	9

rewards

			1
		1	1
		1	
		-10	
		-10	-10
		-10	

- $\pi(s) = \text{“}\uparrow\text{”}, \forall s$
- $\gamma = 0.9$



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

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

$$V_3^\uparrow(6) = \mathbf{R}(6, \uparrow) + 20\% \gamma \mathbf{V}_2^\uparrow(2) + 80\% \gamma \mathbf{V}_2^\uparrow(3)$$

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

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

$(h - 1)$ horizon future values at a next state \mathbf{s}'

sum up future values weighted by the probability of getting to that next state \mathbf{s}'

discounted by γ

Bellman Recursion

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

$$V_1^\uparrow(s) = R(s, \uparrow)$$

0.00	0.00	1.00
0.00	0.00	-10.00
0.00	0.00	0.00

$$V_2^\uparrow(s)$$

0.00	0.00	1.90
0.00	0.00	-9.28
0.00	0.00	-9.00

$$V_3^\uparrow(s)$$

0.00	0.00	2.71
0.00	0.00	-8.63
0.00	0.00	-8.35

$$V_4^\uparrow(s)$$

0.00	0.00	3.44
0.00	0.00	-8.05
0.00	0.00	-7.77

$$V_5^\uparrow(s)$$

0.00	0.00	4.10
0.00	0.00	-7.52
0.00	0.00	-7.24

$$V_6^\uparrow(s)$$

0.00	0.00	4.69
0.00	0.00	-7.05
0.00	0.00	-6.77

$$V_6^\uparrow(6) = R(6, \uparrow) + \gamma[.2V_5^\uparrow(2) + .8 \times V_5^\uparrow(3)]$$

$$-7.048 = -10 + .9[.2 * 0 + 0.8 * 4.10]$$

...

$$V_{61}^\uparrow(s)$$

0.00	0.00	9.98
0.00	0.00	-2.81
0.00	0.00	-2.53

$$V_{62}^\uparrow(s)$$

0.00	0.00	9.99
0.00	0.00	-2.81
0.00	0.00	-2.53

Bellman Recursion

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

approaches infinity

If the horizon h goes to infinity

Bellman Equations

$$V_\infty^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\infty^\pi(s'), \forall s$$

$V_\infty^\uparrow(s)$

$|\mathcal{S}|$ many linear equations, one equation for each state

0.00	0.00	10.00
0.00	0.00	-2.80
0.00	0.00	-2.52

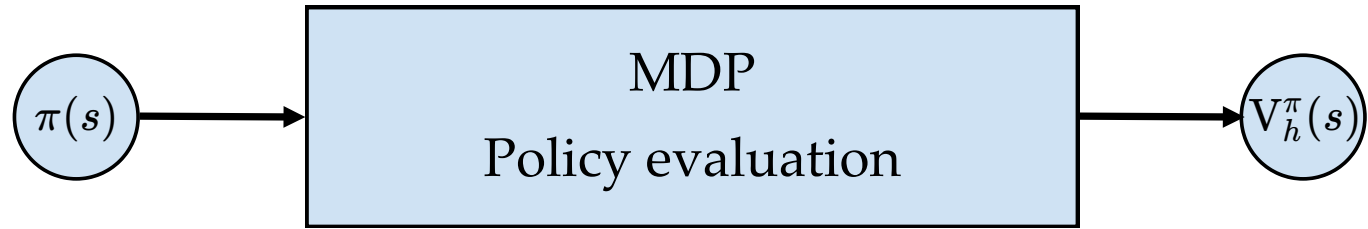
$$10 = V_\infty^\uparrow(3) = R(3, \uparrow) + \gamma[V_\infty^\uparrow(3)] = 1 + .9 \times 10$$

$$-2.8 = V_\infty^\uparrow(6) = R(6, \uparrow) + \gamma[.2V_\infty^\uparrow(2) + .8 \times V_\infty^\uparrow(3)] = -10 + .9[.2 \times 0 + .8 \times 10]$$

$$-2.52 = V_\infty^\uparrow(9) = R(9, \uparrow) + \gamma[V_\infty^\uparrow(6)] = 0 + .9 \times (-2.8)$$

typically $\gamma < 1$ in MDP definition, motivated to make $V_\infty^\pi(s) := \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, \pi]$ finite.

Quick summary



1. By summing h terms:

Recall: For a *given* policy $\pi(s)$, the (state) **value functions**

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

2. By leveraging structure:

finite-horizon Bellman recursions

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

infinite-horizon Bellman equations

$$V_\infty^\pi(s) = \mathbf{R}(s, \pi(s)) + \gamma \sum_{s'} \mathbf{T}(s, \pi(s), s') V_\infty^\pi(s'), \forall s$$

Outline

- Markov Decision Processes Definition, terminologies, and policy
- 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^*

Definition: for a given MDP and a fixed horizon h (possibly infinite), π^* is an optimal policy if $V_h^{\pi^*}(s) = V_h^*(s) \geq V_h^\pi(s)$ for all $s \in \mathcal{S}$ and for all possible policy π .

- An MDP has a unique optimal value $V_h^*(s)$.
- Optimal policy π^* might not be unique (think, e.g. symmetric world).
- For finite h , optimal policy π_h^* depends on how many time steps left.
- When $h \rightarrow \infty$, time no longer matters, i.e., there exists a stationary π^* .
- Under optimal policy, recursion holds too

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

Definition: for a given MDP and a fixed horizon h (possibly infinite), π^* is an optimal policy if $V_h^{\pi^*}(s) = V_h^*(s) \geq V_h^\pi(s)$ for all $s \in \mathcal{S}$ and for all possible policy π .

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

How to search for an optimal policy π^* ?

- One idea: enumerate over all π , do policy evaluation, compare V^π , get $V^*(s)$
- tedious, and even with $V^*(s)$... not super clear how to act

 $V_{61}^*(s)$

8.08	8.98	9.98
7.27	8.08	-1.20
6.54	7.27	6.54

 $V_{62}^*(s)$

8.09	8.99	9.99
7.28	8.09	-1.19
6.55	7.28	6.55

...

 $V_\infty^*(s)$

8.10	9.00	10.00
7.29	8.10	-1.18
6.56	7.29	6.56

Outline

- Markov Decision Processes Definition, terminologies, and policy
- 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 state-action value functions $Q_h^*(s, a)$

$Q_h^*(s, a)$: the expected sum of discounted rewards for

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

$$Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a'), \forall s, a, h$$



recursively finding $Q_h^*(s, a)$

$Q_h^*(s, a)$: the expected sum of discounted rewards for

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

$$Q_0^*(s, a)$$

0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0
0	0	0

$$Q_1^*(s, a) = R(s, a)$$

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Transitions: 2 to 3 (80%), 2 to 5 (20%), 5 to 6 (vertical arrow)

$R(s, a)$

		1
		1
		1
		-10
		-10
		-10



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
	20%	80%
4	5	6
7	8	9

$$Q_1^*(s, a) = R(s, a)$$

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

$$Q_2^*(s, a)$$

				1.9

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

- receive $R(3, \rightarrow)$
- next state $s' = 3$, act **optimally** for the remaining one timestep
 - receive $\max_{a'} Q_1^*(3, a')$

$$\begin{aligned}
 Q_2^*(3, \rightarrow) &= R(3, \rightarrow) + \gamma \max_{a'} Q_1^*(3, a') \\
 &= 1 + .9 \max_{a'} Q_1^*(3, a') \\
 &= 1.9
 \end{aligned}$$



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states (1-9). A dashed arrow points from state 2 to state 3, labeled '80%'. A solid arrow points from state 2 to state 5, labeled '20%'.

$$Q_1^*(s, a) = R(s, a)$$

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

$$Q_2^*(s, a)$$

			1.9	1.9

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

- receive $R(3, \uparrow)$

- next state $s' = 3$, act **optimally** for the remaining one timestep

- receive $\max_{a'} Q_1^*(3, a')$

$$Q_2^*(3, \uparrow) = R(3, \uparrow) + \gamma \max_{a'} Q_1^*(3, a')$$

$$= 1 + .9 \max_{a'} Q_1^*(3, a')$$

$$= 1.9$$



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states (1-9). A dashed arrow points from state 3 to state 2, labeled 20%. A solid arrow points from state 3 to state 6, labeled 80%.

$$Q_1^*(s, a) = R(s, a)$$

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

$$Q_2^*(s, a)$$

				1.9	1.9
				1	1.9

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

- receive $R(3, \leftarrow)$
- next state $s' = 2$, act **optimally** for the remaining one timestep
 - receive $\max_{a'} Q_1^*(2, a')$

$$\begin{aligned}
 Q_2^*(3, \leftarrow) &= R(3, \leftarrow) + \gamma \max_{a'} Q_1^*(2, a') \\
 &= 1 + .9 \max_{a'} Q_1^*(2, a') \\
 &= 1
 \end{aligned}$$



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states (1-9). A dashed arrow points from state 2 to state 5 with a label '20%'. A solid arrow points from state 3 to state 6 with a label '80%'.

$$Q_1^*(s, a) = R(s, a)$$

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

$$Q_2^*(s, a)$$

			1.9	
			1	1.9
				-8

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

- receive $R(3, \downarrow)$

- next state $s' = 6$, act **optimally** for the remaining one timestep

- receive $\max_{a'} Q_1^*(6, a')$

$$Q_2^*(3, \downarrow) = R(3, \downarrow) + \gamma \max_{a'} Q_1^*(6, a')$$

$$= 1 + .9 \max_{a'} Q_1^*(6, a')$$

$$= -8$$



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states (1-9). A dashed arrow points from state 3 to state 2, labeled '20%'. A solid arrow points from state 6 to state 3, labeled '80%'.

$$Q_1^*(s, a) = R(s, a)$$

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

$$Q_2^*(s, a)$$

				1.9	
				1	1.9
				-8	
				-9.28	

- receive $R(6, \uparrow)$
- act **optimally** for one more timestep, at the next state s'

- 20% chance, $s' = 2$, act optimally, receive $\max_{a'} Q_1^*(2, a')$

- 80% chance, $s' = 3$, act optimally, receive $\max_{a'} Q_1^*(3, a')$

Let's consider $Q_2^*(6, \uparrow) = R(6, \uparrow) + \gamma[.2 \max_{a'} Q_1^*(2, a') + .8 \max_{a'} Q_1^*(3, a')]$

$$= -10 + .9[.2 \times 0 + .8 \times 1] = -9.28$$



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states. State 2 is highlighted with a dashed arrow pointing to state 3, labeled '80%'. A solid arrow points from state 5 to state 6, labeled '20%'.

$Q_1^*(s, a)$
= $R(s, a)$

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

$Q_2^*(s, a)$

				1.9	
				1	1.9
				-8	
				-9.28	

$$Q_2^*(6, \uparrow) = R(6, \uparrow) + \gamma[.2 \max_{a'} Q_1^*(2, a') + .8 \max_{a'} Q_1^*(3, a')]$$

in general

$$Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a'), \forall s, a, h$$



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

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

Recall: $\gamma = 0.9$

States and one special transition:

1	2	3
4	5	6
7	8	9

Diagram showing a 3x3 grid of states (1-9). State 2 is connected to state 3 by a solid arrow labeled 80%. State 2 is also connected to state 5 by a dashed arrow labeled 20%. State 5 is connected to state 6 by a solid arrow.

$Q_1^*(s, a)$
= $R(s, a)$

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

$Q_2^*(s, a)$

		1.9
		1 1.9
		-8
		-9.28

what's the optimal action in state 3, with horizon 2, given by $\pi_2^*(3) = ?$

either up or right

in general

$$\pi_h^*(s) = \arg \max_a Q_h^*(s, a), \forall s, h$$

Outline

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

Given the recursion $Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a')$

we can have an infinite horizon equation

$$Q_\infty^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_\infty^*(s', a')$$

Value Iteration

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

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

3. **while** True:

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

5. $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

6. **if** $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$:

7. **return** Q_{new}

8. $Q_{\text{old}} \leftarrow Q_{\text{new}}$

if run this block h times
and break, then the
returns are exactly Q_h^*

$Q_\infty^*(s, a)$



V values vs. Q values

- V is defined over state space; Q is defined over (state, action) space.
- $V_h^*(s)$ can be derived from $Q_h^*(s, a)$:, and vice versa.
- Q^* is easier to read "optimal actions" from.
- We care more about V^π and Q^*

$V_{61}^\uparrow(s)$

0.00	0.00	9.98
0.00	0.00	-2.81
0.00	0.00	-2.53

 $V_{62}^\uparrow(s)$

0.00	0.00	9.99
0.00	0.00	-2.81
0.00	0.00	-2.53

 $V_\infty^\uparrow(s)$

0.00	0.00	10.00
0.00	0.00	-2.80
0.00	0.00	-2.52

 $V_{61}^*(s)$

8.08	8.98	9.98
7.27	8.08	-1.20
6.54	7.27	6.54

 $V_{62}^*(s)$

8.09	8.99	9.99
7.28	8.09	-1.19
6.55	7.28	6.55

 $V_\infty^*(s)$

8.10	9.00	10.00
7.29	8.10	-1.18
6.56	7.29	6.56

$V_h^*(s) = \max_a [Q_h^*(s, a)]$

 $Q_{61}^*(s, a)$

7.27	8.08	9.98			
7.27	8.08	7.27	8.98	9.08	9.98
6.54		7.27		-0.08	
7.27	8.08	-1.20			
6.54	7.27	6.54	-1.08	-2.73	-11.08
5.89		6.54		-4.11	
6.54	7.27	-1.08			
5.89	6.54	5.89	5.89	6.54	5.89
5.89		6.54		5.89	

 $Q_{62}^*(s, a)$

7.28	8.09	9.99			
7.28	8.09	7.28	8.99	9.09	9.99
6.55		7.28		-0.08	
7.28	8.09	-1.19			
6.55	7.28	6.55	-1.08	-2.72	-11.08
5.89		6.55		-4.11	
6.55	7.28	-1.08			
5.89	6.55	5.89	5.89	6.55	5.89
5.89		6.55		5.89	

 $Q_\infty^*(s, a)$

7.29	8.10	10.00			
7.29	8.10	7.29	9.00	9.10	10.00
6.56		7.29		-0.06	
7.29	8.10	-1.18			
6.56	7.29	6.56	-1.06	-2.71	-11.06
5.90		6.56		-4.10	
6.56	7.29	-1.06			
5.90	6.56	5.90	5.90	6.56	5.90
5.90		6.56		5.90	

$\pi_h^*(s) = \arg \max_a [Q_h^*(s, a)]$

Summary

- Markov decision processes (MDP) is 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 , we can calculate $V^\pi(s)$ via
 - the summation over rewards definition
 - Bellman recursion for finite horizon, 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.

<https://forms.gle/DefAyvq8KA9kg37X8>

We'd love to hear
your thoughts.

Thanks!