

Markov Decision Processes

Reinforcement Learning

Roberto Capobianco



SAPIENZA
UNIVERSITÀ DI ROMA

Reinforcement Learning Overview (recap)

— — —

	AI Planning	SL	UL	RL	IL
Optimization	X			X	X
Learns from experience		X	X	X	X
Generalization	X	X	X	X	X
Delayed Consequences	X			X	X
Exploration				X	

- SL = Supervised learning; UL = Unsupervised learning; RL = Reinforcement Learning; IL = Imitation Learning
- Imitation learning assumes input demonstrations of good policies
- IL reduces RL to SL. IL + RL is promising area

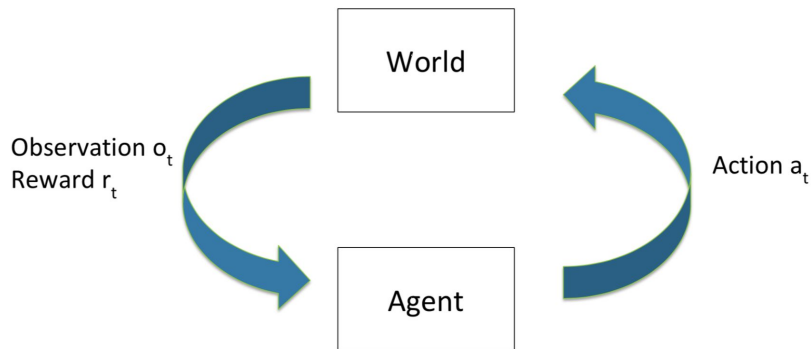
Credits: Emma Brunskill



SAPIENZA
UNIVERSITÀ DI ROMA

Sequential Decision Making

— — —



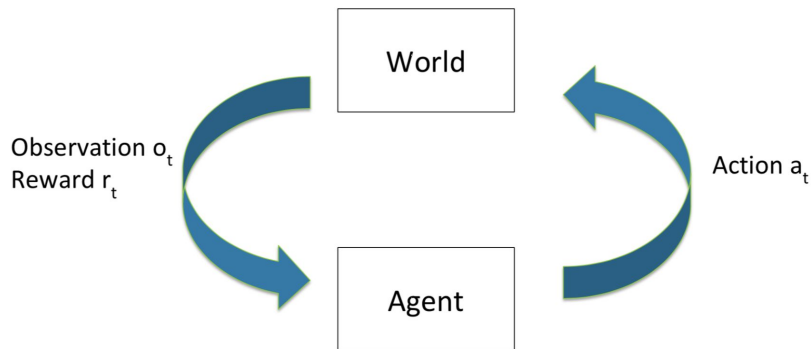
The agent interacts with the environment:

- at discrete timesteps;
- by receiving observations o_t and reward r_t from the environment;
- by taking actions a_t ;



Sequential Decision Making

— — —



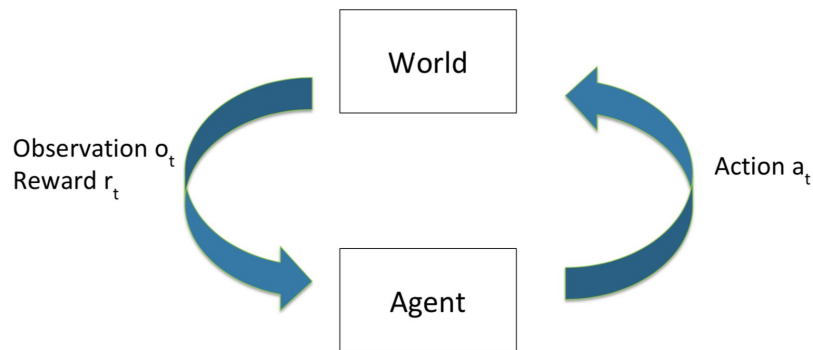
Such discrete interaction generates a trajectory, or history at each timestep t , that is used by the agent to take action:

$$h_t = (o_0, a_0, r_1, o_1, a_1, \dots, r_t, o_t, a_t)$$



Sequential Decision Making

— — —



The state is a function of the history:

$$s_t = f(h_t)$$

and it is typically hidden or unknown



Markov Assumption

A state s_t is Markovian iff future is independent of the past given the present

$$p(s_{t+1} | s_t, a_t) = p(s_{t+1} | h_t, a_t)$$



Markov Assumption

A state s_t is Markovian iff future is independent of the past given the present

$$p(s_{t+1}|s_t, a_t) = p(s_{t+1}|h_t, a_t)$$



Is this problem
Markovian?



Markov Assumption

— — —

- A state can always be made markovian by setting it to be equal to the history

$$s_t = h_t$$

- The best case (used in practice) is: current state corresponds to (or is a sufficient statistic of) latest observation

$$s_t = o_t$$

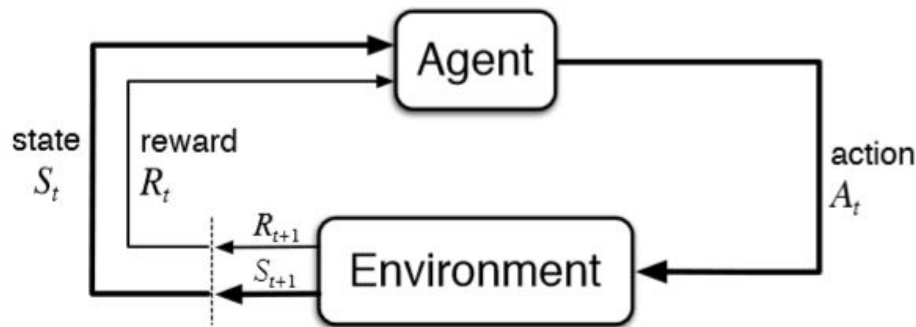
- In this case the state is said to be *fully observable*



Markov Decision Process (MDP)

— — —

- Set of states S
- Set of actions A



Sequential Decision Making under Markov Assumption

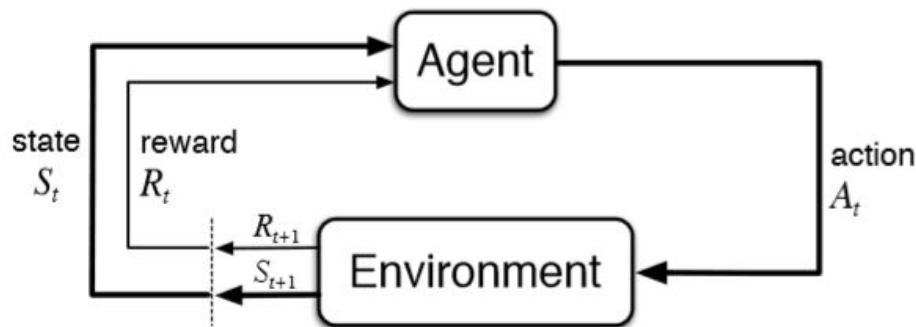
- Markovian transition dynamics
- Full Observability
- The transition dynamics T is (generally) stochastic $p(s_{t+1}|s_t, a_t)$



Markov Decision Process (MDP)

— — —

- Set of states S
- Set of actions A



Alternative notation

Sequential Decision Making under Markov Assumption $s_{t+1} \sim p(\cdot | s_t, a_t)$ or

- Markovian transition dynamics
- Full Observability
- The transition dynamics T is (generally) stochastic $p(s_{t+1} | s_t, a_t)$

$s' \sim p(\cdot | s, a)$



Reward

— — —

A reward r_t is a:

- scalar signal representing a feedback
- indicates how well an agent is doing at step t
- the reward is a function of state and action (often indicated as $R(s,a)$ and sometimes $R(s',a,s)$)
- cost is the inverse of the reward

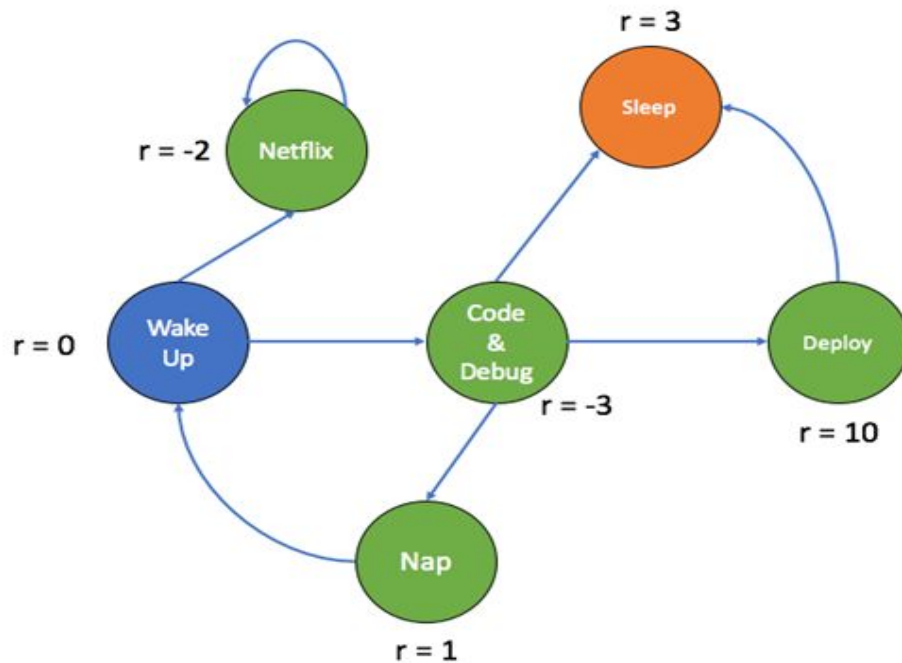
Reward hypothesis: can all goals be achieved through the maximization of a numerical reward?

It's an open question



Deterministic MDP Example

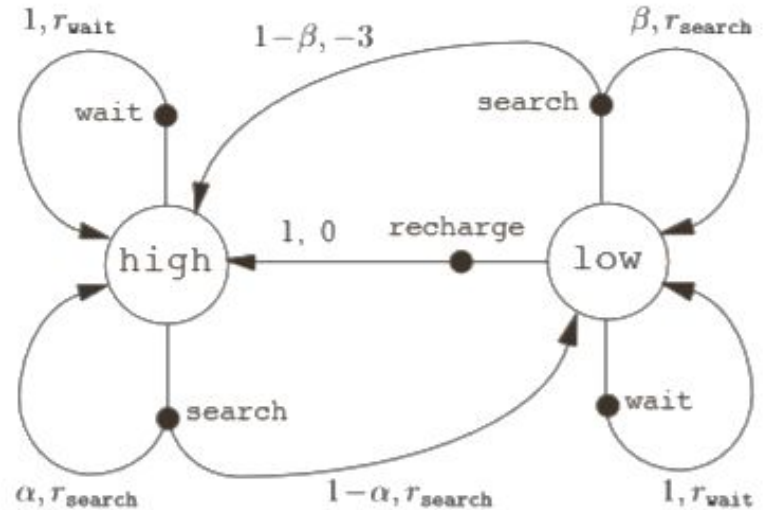
— — —



Stochastic MDP Example

Recycling robot

s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	r_{search}
low	search	high	$1 - \beta$	-3
low	search	low	β	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	$-$
low	wait	high	0	$-$
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	$-$



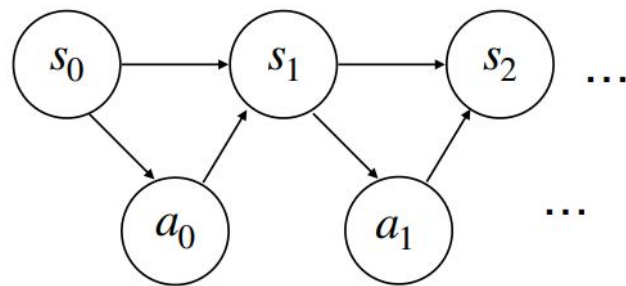
Policy

A policy π :

- is a mapping from (all) states to actions;
- determines how agents select actions;
- can be deterministic ($a = \pi(s)$) or stochastic ($\pi(a|s)$ or $p(a|s)$ or $a \sim \pi(.|s)$)



Trajectory Probability



What's the probability of seeing a trajectory at time t according to π starting at s_0 ?

$$(s_0, a_0, s_1, a_1, \dots, s_t, a_t)$$

$$P^\pi(s_0, a_0, \dots, s_t, a_t) = \pi(a_0 | s_0) p(s_1 | s_0, a_0) \pi(a_1 | s_1) p(s_2 | s_1, a_1) \dots p(s_t | s_{t-1}, a_{t-1}) \pi(a_t | s_t)$$



State Visitation Probability

— — —

What's the probability of visiting state s , a at time t according to π starting at s_0 ?

$$\mathbb{P}_t^\pi(s, a; s_0) = \sum_{a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1}} \mathbb{P}^\pi(s_0, a_0, \dots, s_t=s, a_t=a)$$



Another Example MDP

— — —



- **state:** robot configuration (joint states) and ball position
- **action:** torque on arm and finger joints
- **transition:** stochastic, physics plus noise
- **policy:** mapping from robot state and ball position to torque
- **cost:** magnitude of the torque and distance to the goal



Infinite Horizon Discounted Setting

So far in our MDP we have (S, A, T, R)

Now we add the discount factor γ to reason on the policy's long term effects

- γ is in $[0, 1]$
- $\gamma = 0$ means: I only care about immediate rewards
- $\gamma = 1$ means: Immediate and future rewards are equally important

How so?

Value Function

— — —

- We estimate the goodness of states and actions based on their value
- It's also a measure to compare policies

$$V^{\pi}(s_t) = \mathbb{E}_{\pi}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t] = \mathbb{E}[\sum_{h=0}^{\infty} \gamma^h r_h | s_0 = s_t, a_h = \pi(s_h), s_{h+1} \sim p(\cdot | s_h, a_h)]$$



Value Function/Q-Function

- We estimate the goodness of states and actions based on their value
- It's also a measure to compare policies

$$V^{\pi}(s_t) = \mathbb{E}_{\pi}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t] = \mathbb{E}[\sum_{h=0}^{\infty} \gamma^h r_h | s_0 = s_t, a_h = \pi(s_h), s_{h+1} \sim p(\cdot | s_h, a_h)]$$

$$Q^{\pi}(s_t, a_t) = \mathbb{E}[\sum_{h=0}^{\infty} \gamma^h r_h | (s_0, a_0) = (s_t, a_t), a_{h+1} = \pi(s_h), s_{h+1} \sim p(\cdot | s_h, a_h)]$$



Back to Discount Factor

Setting $\gamma = 1$ for infinite tasks is a bad idea!

Note that $\sum_{h=0}^{\infty} \gamma^h$ is a geometric series and for γ in $[0,1]$ this is equivalent to $1/(1-\gamma)$

So, the value of γ approximately determines how many steps ahead we are considering

E.g., $\gamma=0.99 \rightarrow 99$ timesteps ahead



Bellman Equation

The value of a certain state is expanded in terms of the current reward and the value of the next states according to the policy

$$V^\pi(s_t) = \mathbb{E}_\pi[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t] = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s_t, \pi(s_t))} [V^\pi(s')]$$



Bellman Equation also for Q

The value of a certain state is expanded in terms of the current reward and the value of the next states according to the policy

$$V^{\pi}(s_t) = \mathbb{E}_{\pi}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t] = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [V^{\pi}(s')]$$

$$Q^{\pi}(s_t, a) = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [V^{\pi}(s')]$$



Bellman Equation also for Q

The value of a certain state is expanded in terms of the current reward and the value of the next states according to the policy

r here is function of s and $\pi(s)$

$$V^\pi(s_t) = \mathbb{E}_\pi[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t] = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [V^\pi(s')]$$

$$Q^\pi(s_t, a) = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [V^\pi(s')]$$

r here is function of s and a



Bellman Equation also for Q

The value of a certain state is expanded in terms of the current reward and the value of the next states according to the policy

r here is function of s and $\pi(s)$

$$V^\pi(s_t) = \mathbb{E}_\pi[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t] = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [V^\pi(s')]$$

$$Q^\pi(s_t, a) = r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [V^\pi(s')]$$

r here is function of s and a

As a result $V(s) = Q(s, \pi(s))$



Discounted State-Action Distribution

— — —

$$d^{\pi}_{s_0}(s, a) = (1-\gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s, a; s_0)$$



Discounted State-Action Distribution

— — —

$$d^{\pi}_{s_0}(s, a) = (1-\gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}^{\pi}_h(s, a; s_0)$$

This gives us a probability distribution



Optimal Policy

For infinite horizon MDPs there always exists a deterministic policy π^* such that

$$V^{\pi^*}(s) \geq V^{\pi}(s) \quad \forall s, \pi$$

meaning that π^* dominates all other policies π in each state



Optimal Policy

For infinite horizon MDPs there always exists a deterministic policy π^* such that

$$V^{\pi^*}(s) \geq V^{\pi}(s) \quad \forall s, \pi$$

Alternative notation
 $V^{\pi^*} = V^*$ and $Q^{\pi^*} = Q^*$

meaning that π^* dominates all other policies π in each state



Bellman Optimality

— — —

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]$$



Bellman Optimality

— — —

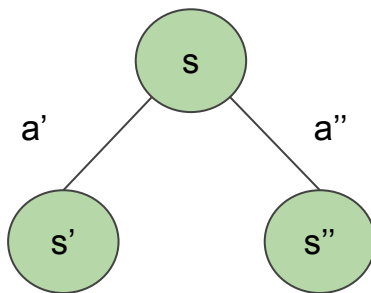
$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]]$$

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



Bellman Optimality Example

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]$$



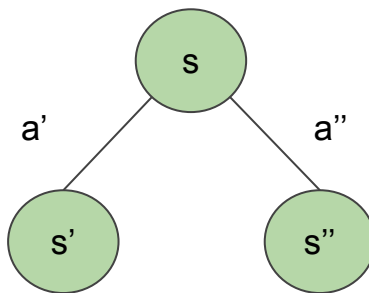
Assume we know V^* at s' and s''



Bellman Optimality Example

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]$$

- Try a' , get $r(s, a')$,
compute
 $Q^*(s, a') = r(s, a') + \gamma V^*(s')$
- Try a'' , get $r(s, a'')$,
compute
 $Q^*(s, a'') = r(s, a'') + \gamma V^*(s'')$



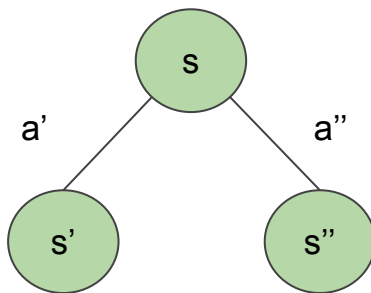
Assume we know V^* at
 s' and s''



Bellman Optimality Example

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]]$$

- Try a' , get $r(s, a')$,
compute
 $Q^*(s, a') = r(s, a') + \gamma V^*(s')$
- Try a'' , get $r(s, a'')$,
compute
 $Q^*(s, a'') = r(s, a'') + \gamma V^*(s'')$



Assume we know V^* at
 s' and s''

$$V^*(s) = \max_{a', a'', \dots} \{ Q^*(s, a'), Q^*(s, a'') \}$$



Bellman Optimality (Theorem 1)

— — —

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]$$

given $\hat{\pi} = \arg\max_a Q^*(s, a)$, we can show $\hat{V}^{\hat{\pi}} = V^*$



Bellman Optimality (Theorem 1)

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]]$$


given $\hat{\pi} = \operatorname{argmax}_a Q^*(s, a)$, we can show $V^{\hat{\pi}} = V^*$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right] \\ &\leq \mathbb{E} [r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots] = V^{\hat{\pi}}(s) \end{aligned}$$



Bellman Optimality (Theorem 1)

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]]$$

given $\hat{\pi} = \operatorname{argmax}_a Q^*(s, a)$, we can show $V^{\hat{\pi}} = V^*$  $V^{\hat{\pi}} \geq V^*$ and $V^* \geq V^{\hat{\pi}}$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right] \\ &\leq \mathbb{E} [r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots] = V^{\hat{\pi}}(s) \end{aligned}$$



Bellman Optimality (Theorem 1)

$$V^*(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V^*(s')]$$

given $\hat{\pi} = \operatorname{argmax}_a Q^*(s, a)$, we can show $\hat{V}^{\hat{\pi}} = V^*$

This implies $\pi^* = \operatorname{argmax}_a Q^*(s, a)$ is an optimal policy



Bellman Optimality (Theorem 2)

For any V , if $V(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V(s')] for all s ,
then $V(s) = V^*(s)$$



Bellman Optimality (Theorem 2)

For any V , if $V(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V(s')] for all s ,
then $V(s) = V^*(s)$$

We need to check if $|V(s) - V^*(s)| = 0$



Bellman Optimality (Theorem 2)

For any V , if $V(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V(s')] for all s ,
then $V(s) = V^*(s)$$

We need to check if

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} |V(s'') - V^*(s'')| \right) \\ &\leq \max_{a_1, a_2, \dots, a_{k-1}} \gamma^k \mathbb{E}_{s_k} |V(s_k) - V^*(s_k)| \end{aligned}$$



Bellman Optimality (Theorem 2)

For any V , if $V(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} V(s')] for all s ,
then $V(s) = V^*(s)$$

$$\begin{aligned} \text{We need to check if } |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} |V(s'') - V^*(s'')| \right) \end{aligned}$$

At infinity, this goes to zero

$$\leq \max_{a_1, a_2, \dots, a_{k-1}} \gamma^k \mathbb{E}_{s_k} |V(s_k) - V^*(s_k)|$$



Bellman Optimality (Theorem 2)

For any V , if $V(s) = \max_a [r(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} V(s')] for all s ,
then $V(s) = V^*(s)$$

This means we can focus on one step at each time (leaving the remaining “problem” to $V(s')$), and any V that satisfies this formula is in fact V^*

