

RL Theory¹: Lecture 1 (Chapter 1)

Shivam Garg, RLAI@UAlberta

27th August 2020

¹based on <https://rltheorybook.github.io/>

Adminstrivia (to set the mood right)

- ▶ Inspiration: Sal Khan and Mark Schmidt.
- ▶ My goal with these lectures is that **all of you** understand most of the material.
- ▶ This partly comes from my own frustration with how inaccessible theory seems to me.
- ▶ So I will go slowly and try to be ~~shameless~~ about it. However, I have this problem where I start picking pace without realizing.
- ▶ Do stop me at any time if something is unclear (be it my accent, a mistake, some skipped steps).
- ▶ And most importantly: "There are 4 kinds of people."

KL



Theory

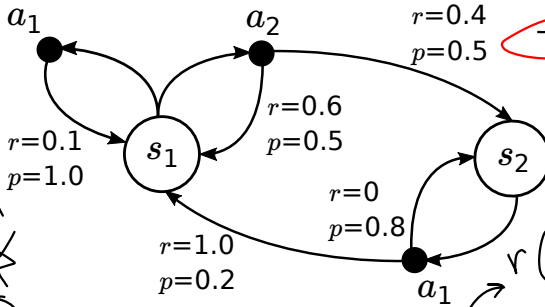
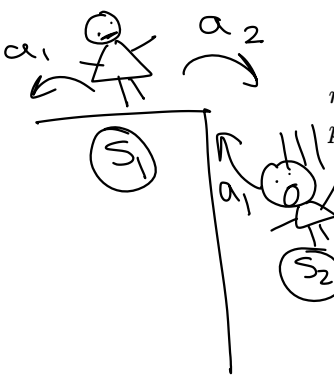


Markov Decision Process

Agent and its Policy

$M(s_1) = 1$ \rightarrow formalize world

$M(s_2) = 0$



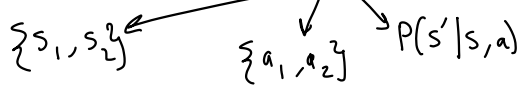
Stick figure $\rightarrow \pi$
 $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$

$\pi(a_1 | s_1) = 1$
 $\pi(a_2 | s_1) = 0$

$\pi(a_1 | s_2) = 1$

$\pi: \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$

Figure: (Rooftop) MDP $M = (\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$.



Agent's Interaction and the Trajectory (\sim stream of experience)

$s_1, a_1, 0.1, s_1, a_1, 0.1, \dots$

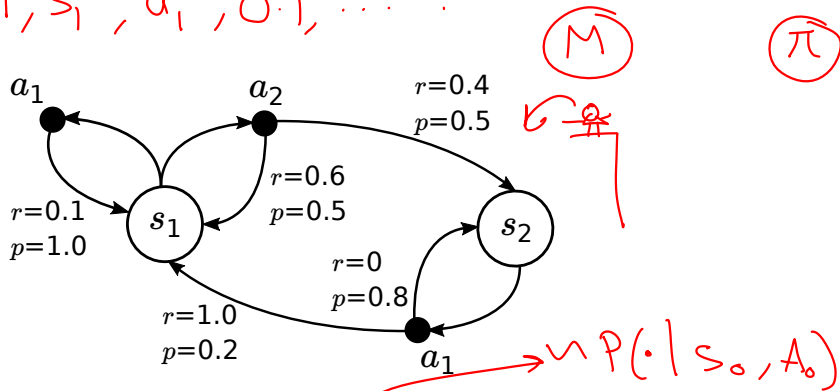


Figure: Trajectory $\tau_t = (S_0, A_0, R_1, S_1, A_1, R_2, S_2, \dots, S_t)$.

$\sim \mathcal{M}(\cdot)$ $\sim \pi(\cdot | s_0)$ $r(S_0, A_0)$

Goal of the Agent ^{safe} π

$$R \in [0, 1]$$

$$1 + \gamma + \gamma^2 + \dots = \frac{1}{1 - \gamma}$$

- ▶ Agent interacts with the environment to generate a trajectory

$$\tau_t = (S_0, A_0, R_1, S_1, A_1, R_2, S_2, \dots, S_t).$$

$$\|Q^* - Q\|_\infty$$

$$\leq \frac{1}{(1 - \gamma)}$$

- ▶ Then agent learns to optimize its policy π to maximize the return (in expectation)

$$G_0 = \frac{1}{1 - \gamma} (R_1 + \gamma R_2 + \gamma^2 R_3 + \dots)$$

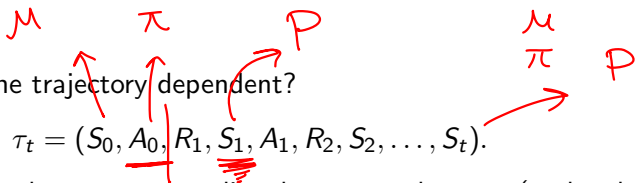
safe π

- ▶ Example: For the trajectory $\tau = (s_1, a_1, 0.1, s_1, a_1, 0.1, \dots)$ and $\gamma = 0.9$,

$$G_0 = (1 - \gamma) [0.1 + \gamma \cdot 0.1 + \gamma^2 \cdot 0.1 + \dots]$$

$$= [0.1]$$

State Value function



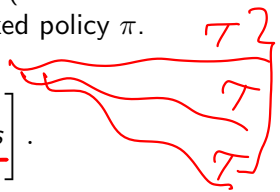
- ▶ On what factors is the trajectory dependent?

$$\tau_t = (S_0, A_0, R_1, S_1, A_1, R_2, S_2, \dots, S_t)$$

- ▶ Value function allows the agent to predict the expected return (under the environment's transition dynamics) from a given state for a fixed policy π .
- ▶ State Value function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$

$$V^\pi(s) = (1 - \gamma) \mathbb{E}_{P, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) \mid S_0 = s \right]$$

Handwritten annotations: $V^\pi(s)$ is circled, $(1 - \gamma)$ is boxed, and $\mathbb{E}_{P, \pi}$ is underlined. A red arrow points from π to the expectation operator. A bracket under the summation indicates it is a summation.



- ▶ Expectation \rightarrow summation notation:

$$V^\pi(s) = \sum_{a_0} \pi(a_0 | s) \left[(1 - \gamma) r(s, a_0) + \sum_{s_1} P(s_1 | s, a_0) \left[\sum_{a_1} \pi(a_1 | s_1) \dots \right] \right]$$

Handwritten equation showing the expansion of the expectation operator into a summation over actions a_0 and states s_1 .

Policy Evaluation: V^π for Rooftop MDP

$$(1-\gamma) \mathbb{E}_{\pi, P} \left[\sum \gamma^t R \mid S_0 = s \right]$$

- Consider the safe policy: $\pi(A|s_1) = \begin{cases} 1 & A = a_1, \\ 0 & A = a_2. \end{cases}$

and $\pi(a_2|s_2) = 1$.

- Then $V^\pi(s_1) = 0.1$

- What about $V^\pi(s_2) = ?$

π

τ

0.1

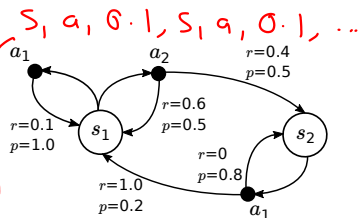


Figure: (Rooftop) MDP.

$P(\tau)$	τ
0.2	$s_2, a_1, r, s_1, a, 0.1, \dots$
0.8×0.2	$s_2, a_1, 0, s_2, a_1, 0.1, s_1, \dots$
$0.8^2 \times 0.2$	s_2, s_2, s_1
\vdots	\vdots

(Policy evaluation is tedious! We'll later describe how to do it iteratively.)

Goal of RL Agent (again!)

s_0

$$\pi^* = \max_{\pi \in \Pi} V^\pi(s) \quad \text{for given state } s$$

$$\pi_1 \quad \text{---} \quad V^{\pi_1}(s)$$

$$\pi_2 \quad \text{---} \quad V^{\pi_2}(s)$$

$$\pi^* \quad \text{---}$$

State–Action Value Function

S

$a \sim \pi$

S

$$Q^\pi(s, a) = (1 - \gamma) \mathbb{E}_{P, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid S_0 = s, A_0 = a \right].$$

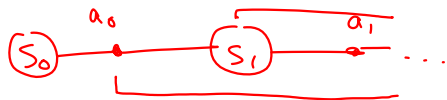
given
 a

S

—•—

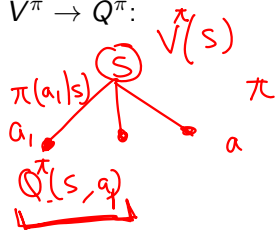
...
 π

Bellman (Consistency) Equations



$V^\pi(s) = Q^\pi(s, \pi(s))$ for deterministic policy π .

► $V^\pi \rightarrow Q^\pi$:



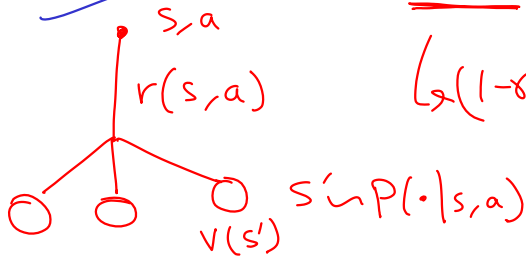
$$V^\pi(s) = \sum_a \pi(a | s) Q^\pi(s, a)$$

$[0, 1]$

normalized

► $Q^\pi \rightarrow V^\pi$:

$$Q^\pi(s, a) = (1 - \gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^\pi(s')].$$



$$\Leftrightarrow (1 - \gamma)r(s, a) + \gamma \sum_{s'} p(s' | s, a) V^\pi(s')$$

Reminder about Matrix-Vector Multiplication

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \times 1 + 2 \times 2 \\ 3 \times 1 + 4 \times 2 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 3 \end{bmatrix} \times 1 + \begin{bmatrix} 2 \\ 4 \end{bmatrix} \times 2 \quad \checkmark$$

$$\begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0$$

Vector Notation for r , V , and Q

- ▶ The reward vector $r \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$:

$$r = \begin{bmatrix} r(s_1, a_1) \\ r(s_1, a_2) \\ r(s_2, a_1) \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.5 \\ 0.2 \end{bmatrix} \quad Q(s, a)$$

- ▶ Value function vectors $V^\pi \in \mathbb{R}^{|\mathcal{S}|}$ and $Q^\pi \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$:

$$V^\pi = \begin{bmatrix} V^\pi(s_1) \\ V^\pi(s_2) \end{bmatrix}; \quad Q^\pi = \begin{bmatrix} Q^\pi(s_1, a_1) \\ Q^\pi(s_1, a_2) \\ Q^\pi(s_2, a_1) \end{bmatrix}.$$

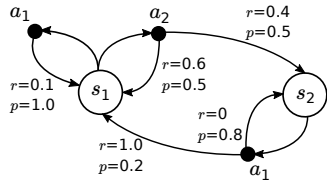


Figure: (Rooftop) MDP.

Vector Notation for P and P^π

P P^π

- $P \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}|}$, with entry $P_{(s,a),s'} = P(s'|s, a)$.

$$(s,a) \begin{bmatrix} P(s_1|s_1,a_1) & P(s_2|s_1,a_1) \\ P(s_1|s_1,a_2) & P(s_2|s_1,a_2) \\ P(s_1|s_2,a_1) & P(s_2|s_2,a_1) \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0.5 & 0.5 \\ 0.2 & 0.8 \end{bmatrix}$$

- $P^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}| \times |\mathcal{A}|}$, with entry $P_{(s,a),(s',a')} = P(s'|s, a) \cdot \pi(a'|s')$.

$$\begin{bmatrix} s_1 a_1 | s_1 a_1 & s_1 a_2 | s_1 a_1 & s_2 a_1 | s_1 a_1 \\ \vdots & \vdots & \vdots \\ 3 \times 3 & & \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

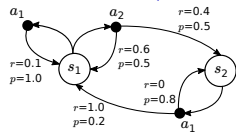
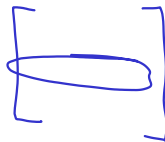
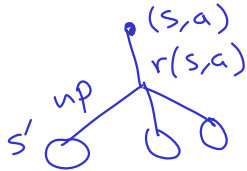


Figure: (**Rooftop**) MDP.

Vector Notation for $Q^\pi \rightarrow V^\pi$ Bellman Equation



$$Q^\pi = (1 - \gamma)r + \gamma PV^\pi$$

Annotations: 3×1 (pointing to Q^π), 2×1 (pointing to r), 3×2 (pointing to PV^π), and 3×1 (pointing to r).

$$\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$Q^\pi(s, a) = (1 - \gamma)r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^\pi(s')$$

$$\begin{bmatrix} Q^\pi(s_1, a_1) \\ Q^\pi(s_1, a_2) \\ Q^\pi(s_2, a_1) \end{bmatrix}$$

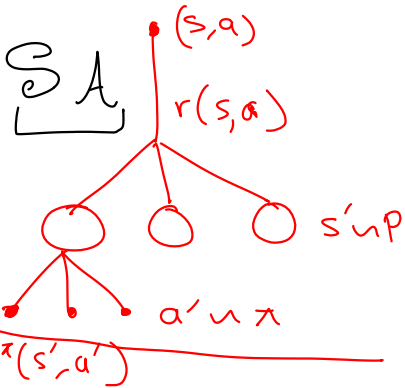
$$= (1 - \gamma) \begin{bmatrix} r(s_1, a_1) \\ \cdot \\ \cdot \end{bmatrix} + \gamma$$

$$\begin{bmatrix} p(s_1|s_1, a_1) & p(s_2|s_1, a_1) \\ p(s_1|s_1, a_2) & p(s_2|s_1, a_2) \\ p(s_1|s_2, a_1) & p(s_2|s_2, a_1) \end{bmatrix} \begin{bmatrix} V(s_1) \\ V(s_2) \end{bmatrix}$$

Vector Notation for $Q^\pi \rightarrow Q^\pi$ Bellman Equation

$$\begin{bmatrix} P(s_1, a_1 | s_1, a_1) \\ P(s_1, a_2 | s_1, a_1) \\ P(s_2, a_1 | s_1, a_1) \end{bmatrix}$$

$$Q^\pi = (1 - \gamma)r + \gamma P^\pi Q^\pi$$



$$Q^\pi(s, a) = (1 - \gamma)r(s, a) +$$

$$\gamma \sum_{s'} P(s' | s, a) \sum_{a'} \pi(a' | s') Q^\pi(s', a')$$

$$= (1 - \gamma)r(s, a) + \gamma \sum_{s'} \sum_{a'} P(s' | s, a) \pi(a' | s') Q(s', a')$$

$$Q^\pi(s, a) = (1 - \gamma)r(s, a) + \gamma \sum_{s', a'} P^\pi(s', a' | s, a) Q^\pi(s', a')$$

Exact Solution for Q^π

$$\hookrightarrow Q^\pi = (1 - \gamma)r + \gamma P^\pi Q^\pi.$$

$$Q^\pi - \gamma P^\pi Q^\pi = (1 - \gamma)r$$

$$(I - \gamma P^\pi) Q^\pi = (1 - \gamma)r$$

$$Q^\pi = (I - \gamma P^\pi)^{-1} (1 - \gamma)r$$



Solving Exactly for Q^π (Rooftop MDP with Safe Policy)

- ▶ Calculation:

$$\begin{aligned}
 Q^\pi &= (1 - \gamma) \overbrace{(I - \gamma P^\pi)^{-1}}^{P^\pi} r \\
 &= (1 - 0.9) \left(\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix} \right)^{-1} \begin{bmatrix} 0.1 \\ 0.5 \\ 0.2 \end{bmatrix} \\
 &= \begin{bmatrix} Q^\pi(s_1, a_1) \\ Q^\pi(s_1, a_2) \\ Q^\pi(s_2, a_1) \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1445 \\ 0.11 \end{bmatrix}
 \end{aligned}$$

- ▶ Define $\pi_Q(s) := \arg \max_{a \in \mathcal{A}} Q(s, a)$.
- ▶ Define $V_Q(s) := \max_{a \in \mathcal{A}} Q(s, a)$.

$$V_Q(s_1) = 0.1445$$

$$V_Q(s_2) = 0.11$$

$$V_Q = (1 - \gamma) \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} + \gamma P \begin{bmatrix} 0.1445 \\ 0.11 \end{bmatrix}$$

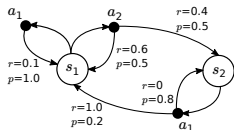
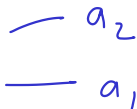


Figure: (Rooftop) MDP.



$I - \gamma P^\pi$ is Invertible

$r \in [0, 1)$

▶ Invertibility = Full Rank

$$\begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \underbrace{A}_A a + \underbrace{B}_B b + \underbrace{C}_C c \neq \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$$

$a, b, c \neq 0$

$$Xx \neq 0 \quad \forall x \neq 0$$

▶ ∞ -norm and triangle inequality

$$\|X\|_\infty = \max_x |x|$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \rightarrow 3$$

$$\|A\|_\infty + \|B\|_\infty \geq \|A+B\|_\infty$$

$$\|A-B\| + \|B\| \geq \|A\|$$

$$\|A-B\| \geq \|A\| - \|B\|$$



max |x|

$(I - \gamma P^\pi)$ is Invertible (Proof)

$$\boxed{x \neq 0} \quad \left\| \underbrace{(I - \gamma P^\pi)x}_{\text{4}} \right\|_\infty = \|x - \gamma P^\pi x\|_\infty$$

$$\geq \|x\|_\infty - \gamma \|P^\pi x\|_\infty$$

$$\geq \|x\|_\infty - \gamma \|x\|_\infty = (1 - \gamma) \|x\|_\infty$$

$$> 0$$

$\gamma \rightarrow 0$

$$+ \|P^\pi x\|_\infty \leq \|x\|_\infty$$

$$- \|P^\pi x\|_\infty \geq -\|x\|_\infty \quad x^* = \max_i x_i = \|x\|_\infty$$

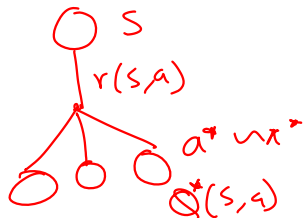
$$\leq p_1 x^* + p_2 x^* + p_3 x^*$$

$$\sum p_i = 1 \quad \uparrow = x^*$$

$$\begin{bmatrix} p_1 & p_2 & p_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} p_1 x_1 + p_2 x_2 + p_3 x_3 \end{bmatrix} \leq \begin{bmatrix} \max x_i \end{bmatrix}$$

Bellman Optimality Equations

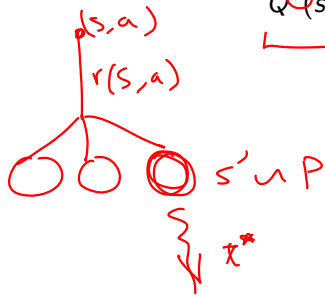
► $V^* \rightarrow Q^*$



$$V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a).$$

$$V^*(s) = Q^*(s, \pi^*(s))$$

► $Q^* \rightarrow V^*$



$$Q^*(s, a) = (1 - \gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')].$$

Bellman Optimality Operator \mathcal{T}

► $\mathcal{T} : \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|} \rightarrow \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$

$Q^*(s_1, a_1)$
 $Q^*(s_1, a_2)$
 $Q^*(s_2, a_1)$

$Q^*(s_1, \pi^*(s_1))$
 \downarrow
 $V^*(s_1)$

$$\mathcal{T}Q := (1 - \gamma)r + \gamma PV_Q.$$

$$\begin{aligned} \mathcal{T}Q(s, a) &= (1 - \gamma)r(s, a) + \gamma \sum_{s'} p(s' | s, a) V_Q(s') \\ &= (1 - \gamma)r(s, a) + \gamma \sum_{s'} p(s' | s, a) \max_{a'} Q(s', a') \end{aligned}$$

$$\mathcal{T}Q(s, a) = (1 - \gamma)r(s, a) + \max_{a'} \sum_{s'} p(s' | s, a) Q(s', a')$$

► Therefore, the $Q^* \rightarrow V^*$ equation can be written as

$$Q^*(s, a) = (1 - \gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')]$$

\Rightarrow

$$Q^* = \mathcal{T}Q^*.$$

Fixed point

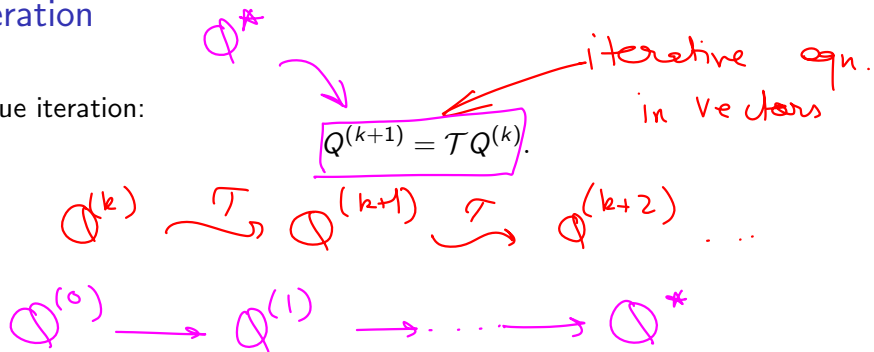
Q^*

$$V_{Q^*} = V^*$$

$$(1 - \gamma)r(s, a) + \gamma \sum_{s'} p(s' | s, a) V_{Q^*}(s')$$

Q-Value Iteration

- ▶ Q-value iteration:



- ▶ What are iterative equations?

- ▶ E.g.: $x_{k+1} = 4x_k + 3$. (divergent for $x_0 = 0$)

$$0, 3, 15, \dots \rightarrow$$

- ▶ E.g.: $y_{k+1} = 0.1y_k + 1$ (convergent for $y_0 = 0$)

$$0, 1, 1.1, 1.11, \dots$$

$$\uparrow \\ 0.111111$$

$$1 + 0.1 \times 0.111111 = 0.111111 + 1 = 1.111111$$

Q-Value Iteration Converges!

$$V^* - V^k \leq 0.1 \quad \checkmark$$

Theorem (Q-Value Iteration Convergence)

- ▶ Set $Q^{(0)} = 0$.
- ▶ Obtain $Q^{(k+1)} = \mathcal{T}Q^{(k)}$ for $k = 0, 1, 2, \dots$
- ▶ Let $\pi^{(k)} = \pi_{Q^{(k)}}$.

$\pi^{(k)}$ is the greedy wrt. $Q^{(k)}$

- ▶ Then for $k > \frac{1}{1-\gamma} \log\left(\frac{2}{\epsilon(1-\gamma)}\right)$,

$$V^* > V^{\pi^{(k)}}$$

$$V^{\pi^{(k)}} \geq V^* - \epsilon \mathbb{1} \equiv$$

$$\boxed{V^*(s) - V^{\pi^{(k)}}(s) \leq \epsilon \quad \forall s}$$

$$V^* - V^{\pi^{(k)}} \leq \epsilon \quad \Downarrow$$

$$O(\log(1/\epsilon))$$

$$\begin{aligned} \epsilon &= 0.1 \\ k &= \log(10) \end{aligned}$$

$$\begin{aligned} \epsilon &= 0.01 \\ k &= \log 100 \end{aligned}$$

Proof: Q-Value Iteration Convergence

- ▶ To prove: For $k > \frac{1}{1-\gamma} \log \left(\frac{2}{\epsilon(1-\gamma)} \right)$, $V^{\pi^{(k)}} \geq V^* - \epsilon \mathbb{1}$ holds.

- ▶ Bellman Optimality Operator \mathcal{T} is a Contraction

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_{\infty} \leq \gamma \|Q - Q'\|_{\infty}.$$

- ▶ Bounds on $\|Q^{(k)} - Q^*\|_{\infty}$

$$\|Q^{(k)} - Q^*\|_{\infty} \leq e^{-(1-\gamma)k}.$$

- ▶ Bounding the Suboptimality of π_Q

$$V^{\pi_Q} \geq V^* - \frac{2}{1-\gamma} \|Q - Q^*\|_{\infty} \mathbb{1}.$$

(Part 1): Bellman Optimality Operator \mathcal{T} is a Contraction

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty.$$

(Part 2): Bounds on $\|Q^{(k)} - Q^*\|_\infty$

$$\|Q^{(k)} - Q^*\|_\infty \leq e^{-(1-\gamma)k}.$$

(Part 3): Bounding the Suboptimality of π_Q

$$V^{\pi_Q} \geq V^* - \frac{2}{1-\gamma} \|Q - Q^*\|_{\infty} \mathbb{1}.$$

(Final Part): Proof of Q -Value Iteration Convergence

Summary

- ▶ What is MDP?
- ▶ What is an agent?
- ▶ The goal of RL.
- ▶ Value iterations.
- ▶ Next time (play with Rooftop MDP, existence of an optimal stationary and deterministic policy, policy iteration)