

Basics of MDPs and Planning in MDPs

Wen Sun

CS 6789: Foundations of Reinforcement Learning

Announcements

Office hours (mine and DAs) are coming

HW0 is due Feb 1st

Recap: RL versus Supervised Learning

Q: How to solve Supervised Learning via RL?

Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$\text{Q function } Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Theorem 1: Bellman Optimality (Q-version)

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[\max_{a' \in A} Q^*(s', a') \right]$$

Main Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$, How to find π^* (stationary & deterministic)

Outline

1. Bellman optimality — property of V^*
2. Optimal planning: Value Iteration

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$|V(s) - V^*(s)| = \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right|$$

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

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

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

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

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

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

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

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

Bellman Optimality for Q^*

What about Q^* ?

Bellman Optimality for Q^*

What about Q^* ?

We should have:

For any $Q : S \times A \rightarrow \mathbb{R}$, if $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q(s', a')$
for all s , then $Q(s, a) = Q^*(s, a), \forall s, a$

Outline

1. Bellman optimality — property of V^*

2. Optimal planning: Value Iteration

Define Bellman Operator \mathcal{T} :

Given a function $f : S \times A \mapsto \mathbb{R}$,

$$\mathcal{T}f : S \times A \mapsto \mathbb{R},$$

$$(\mathcal{T}f)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} f(s', a'), \forall s, a \in S \times A$$

Value Iteration Algorithm:

1. Initialization: $Q^0 : \|Q^0\|_\infty \in (0, \frac{1}{1-\gamma})$
2. Iterate until convergence: $Q^{t+1} = \mathcal{T} Q^t$

Intuition:

Via Bellman optimality theorem:

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

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T}f$

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^* = \ell(x^*)$

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T} f$

Consider the simple problem: finding fixed point solution $x^* = \ell(x^*)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots,$$

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^* = \ell(x^*)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots,$$

$$|x_t - x^*| =$$

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^* = \ell(x^*)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots,$$

$$|x_t - x^*| = |\ell(x_{t-1}) - \ell(x^*)|$$

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^* = \ell(x^*)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots,$$

$$|x_t - x^*| = |\ell(x_{t-1}) - \ell(x^*)| \leq L|x_{t-1} - x^*|$$

Intuition:

Via Bellman optimality theorem:

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

i.e., Q^* is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^* = \ell(x^*)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots,$$

$$|x_t - x^*| = |\ell(x_{t-1}) - \ell(x^*)| \leq L|x_{t-1} - x^*|$$

If $L < 1$ (i.e., contraction), then it converges exponentially fast

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

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

Proof:

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

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

Proof:

$$|\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| = \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right|$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

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

Proof:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \end{aligned}$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

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

Proof:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \end{aligned}$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

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

Proof:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \\ &\leq \gamma \max_{s'} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \end{aligned}$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

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

Proof:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \\ &\leq \gamma \max_{s'} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| = \gamma \|Q - Q'\|_\infty \end{aligned}$$

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Proof:

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Proof:

$$\|Q^{t+1} - Q^*\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^*\|_\infty \leq \gamma \|Q^t - Q^*\|_\infty$$

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Proof:

$$\|Q^{t+1} - Q^*\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^*\|_\infty \leq \gamma \|Q^t - Q^*\|_\infty$$

$$\dots \leq \gamma^{t+1} \|\widehat{Q}^0 - Q^*\|_\infty$$

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$V^{\pi^t}(s) - V^*(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s))$$

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \end{aligned}$$

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \end{aligned}$$

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \end{aligned}$$

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) - 2\gamma^t \|Q^0 - Q^*\|_\infty \end{aligned}$$

Final Quality of the Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) - 2\gamma^t \|Q^0 - Q^*\|_\infty \quad \dots \text{Recursion} \end{aligned}$$

Summary for today

Planning algorithm (no learning so far):

VI: fixed point iteration $Q^{t+1} = \mathcal{T} Q^t$

1. Bellman operator is a contraction map
2. $\|Q^t - Q^*\|_\infty$ being small implies V^{π^t} & V^* are close