

CS 6789: Foundations of Reinforcement Learning

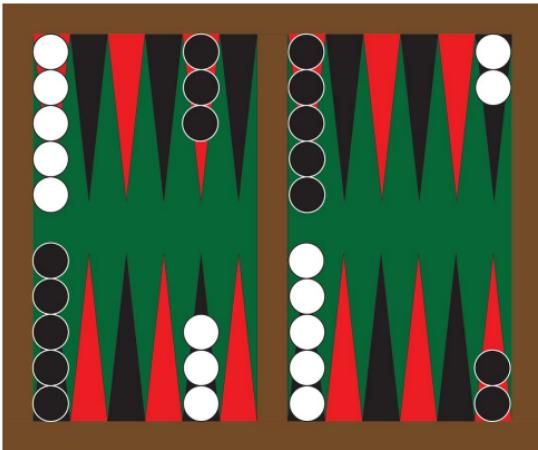
Sham Kakade (UW & MSR) & Wen Sun

TA: Jonathan Chang

<https://wensun.github.io/CS6789.html>

Fall 2020

Progress of RL in Practice



TD GAMMON [Tesauro 95]



[AlphaZero, Silver et.al, 17]



[OpenAI Five, 18]

This course focuses on **RL Theory**

We care about sample complexity:

Total number of environment interactions
needed to learn a high quality policy (i.e., achieve the task)

Four main themes we will cover in this course:

1. Exploration strategies (not just random)
2. Policy Optimization (gradient descent)
3. Control (LQR and nonlinear control)
4. Imitation Learning (i.e., learning from demonstrations)

Logistics

Four (HW0-HW3) assignments (total 55%) and Course Project (45%)

(HW0 10%, HW1-3 15% each)

HW0 is out today and due in one week

Prerequisites (HW0)

Deep understanding of Machine Learning, Optimization, Statistics

ML: sample complexity analysis for supervised learning (PAC)

Opt: Convex (linear) optimization, e.g., gradient decent for convex functions

Stats: basics of concentration (e.g., Hoeffding's), tricks such as union bound

Prerequisites (HW0)

Deep understanding of Machine Learning, Optimization, Statistics

ML: sample complexity analysis for supervised learning (PAC)

Opt: Convex (linear) optimization, e.g., gradient decent for convex functions

Stats: basics of concentration (e.g., Hoeffding's), tricks such as union bound

Check out HW0 asap!

Course projects (45%)

- Team work: size 3
- Midterm report (5%), Final presentation (15%), and Final report (25%)
- Basics: **survey** of a set of similar RL theory papers. Reproduce analysis and provide a coherent story
- Advanced: **identify** extensions of existing RL papers, **formulate** theory questions, and **provide** proofs

Course Notes: Reinforcement Learning Theory & Algorithms

- Book website: <https://rltheorybook.github.io/>
- Many lectures will correspond to chapters in Version 2.
 - V2 (in progress): will use unnormalized rewards.
 - V1: feel free to read ahead.
Be aware V1 uses normalized values functions
- Please do the assigned readings, especially if you want a mastery of the material.
- Please let us know if you find typos/errors in the book!
We appreciate it!

Basics of Markov Decision Processes

Supervised Learning

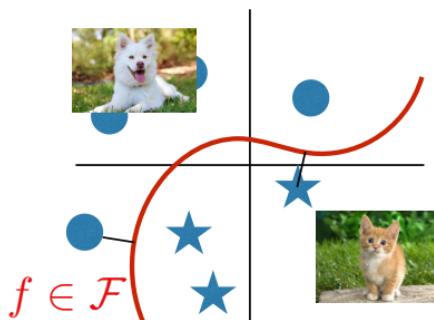
Supervised Learning

Given i.i.d examples at training:



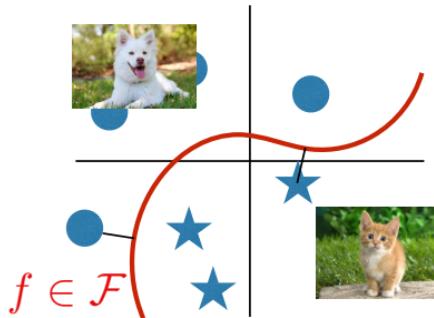
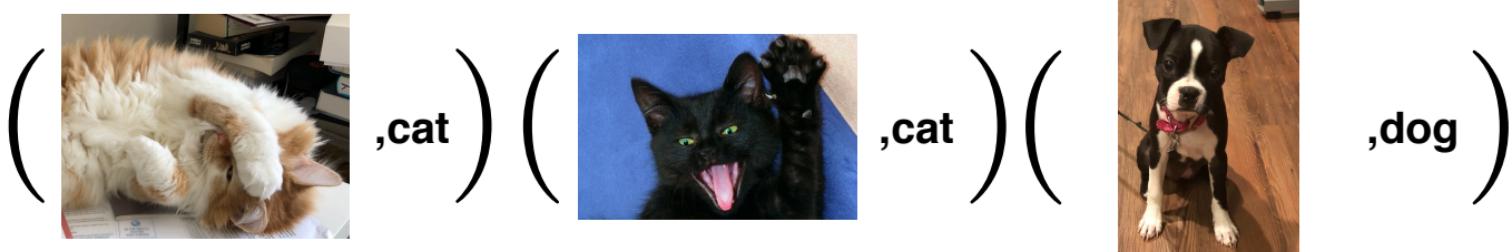
Supervised Learning

Given i.i.d examples at training:

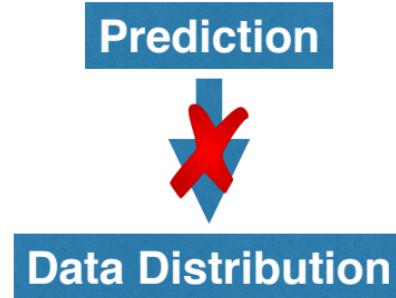


Supervised Learning

Given i.i.d examples at training:



Passive:



Supervised Learning

Given i.i.d examples at training:

Cersei



,cat

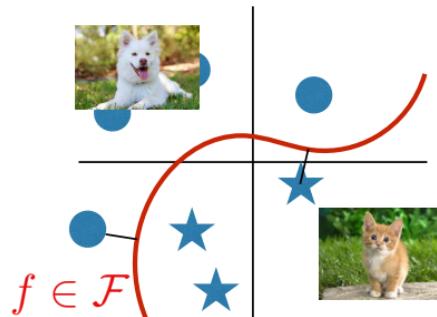


,cat

Tasha



,dog



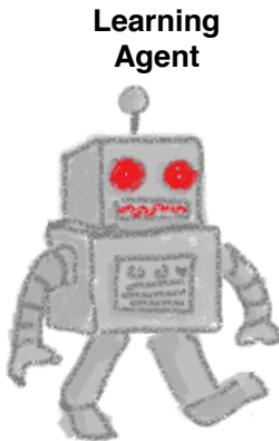
Passive:

Prediction



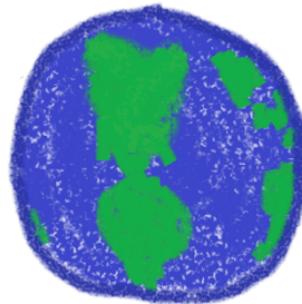
Data Distribution

Markov Decision Process

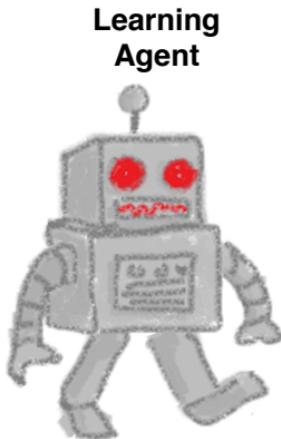


$\pi(s) \rightarrow a$
Policy: determine **action** based on **state**

Environment



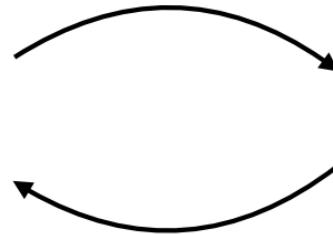
Markov Decision Process



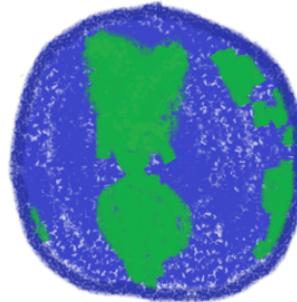
Learning
Agent

$$\pi(s) \rightarrow a$$

Policy: determine **action** based on **state**



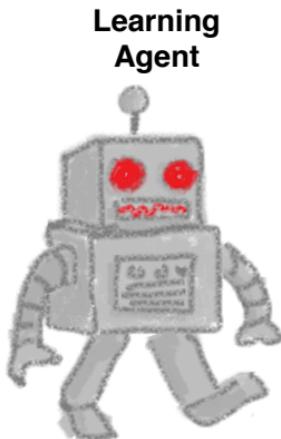
Environment



Send **reward** and **next state** from a
Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

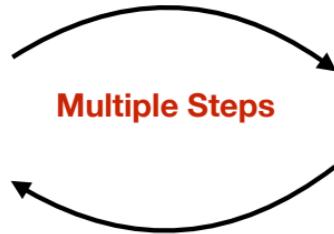
Markov Decision Process



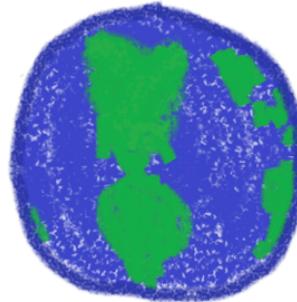
Learning
Agent

$$\pi(s) \rightarrow a$$

Policy: determine **action** based on **state**



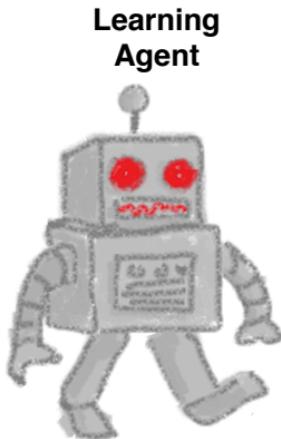
Environment



Send **reward** and **next state** from a
Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

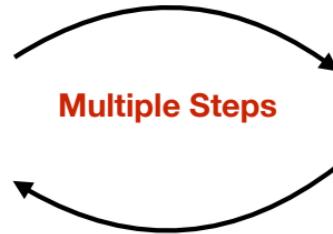
Markov Decision Process



Learning
Agent

$$\pi(s) \rightarrow a$$

Policy: determine **action** based on **state**



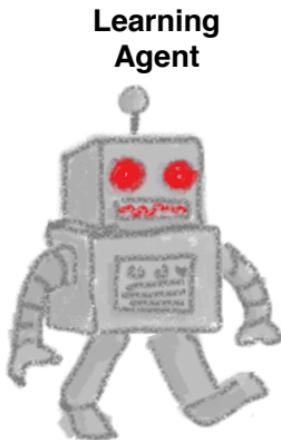
Environment



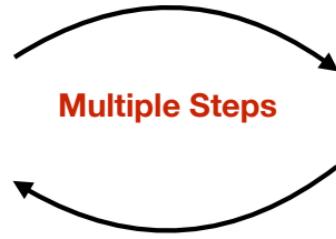
Send **reward** and **next state** from a
Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

Markov Decision Process



$\pi(s) \rightarrow a$
Policy: determine **action** based on **state**



Send **reward** and **next state** from a
Markovian transition dynamics

Environment



$$r(s, a), s' \sim P(\cdot | s, a)$$

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning					
Reinforcement Learning					

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓				
Reinforcement Learning	✓				

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓			

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓	✓		

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓	✓	✓	

	Learn from Experience	Generalize	Interactive	Exploration	Credit assignment
Supervised Learning	✓	✓			
Reinforcement Learning	✓	✓	✓	✓	✓

Infinite horizon Discounted Setting

statespace
Action space

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

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

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

Infinite horizon Discounted Setting

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

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

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

$$\begin{array}{c} \pi(s) \leftarrow \Delta(A) \\ a \sim \pi(s) \\ a \sim \pi(\cdot | s) \end{array}$$

$s_0 \xrightarrow{a_0 \sim \pi(s_0)} s_1 \sim P(\cdot | s_0, a_0)$

$a_1 \sim \pi(\cdot | s_1)$

$s_2 \sim P(\cdot | s_1, a_1)$

\vdots

$r(s_0, a_0)$

$r(s_1, a_1)$

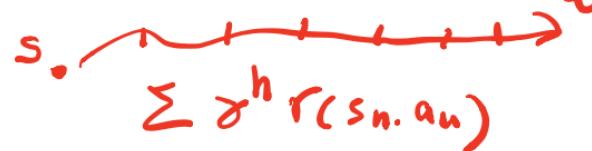
Infinite horizon Discounted Setting

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

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

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

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



Infinite horizon Discounted Setting

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

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

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

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

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



Bellman Equation:

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

Bellman Equation:

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

$$\underbrace{V^\pi(s)}_{\Delta} = \mathbb{E}_{a \sim \pi(s)} \left[\underbrace{r(s, a)}_{\Delta} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^\pi(s') \right]$$



$$V^\pi(s) = \mathbb{E} \left[\underbrace{r(s_0, a_0) + \gamma r(s_1, a_1) + \dots}_{\Delta} \right] \quad V^\pi(s')$$

$$= \mathbb{E}_{a_0 \sim \pi(s_0)} \left[\underbrace{r(s_0, a_0)}_{\Delta} + \gamma \mathbb{E}_{\substack{s_1 \sim P(\cdot | s_0, a_0)}} \left[\underbrace{r(s_1, a_1) + \gamma r(s_2, a_2) + \dots}_{\Delta} \right] \right] = V^\pi(s_0)$$

$$= \mathbb{E}_{a_0 \sim \pi(s_0)} \left[r(s_0, a_0) + \gamma \mathbb{E}_{\substack{s_1, a_1 \sim P(\cdot | s_0, a_0)}} \left[\underbrace{r(s_1, a_1) + \gamma r(s_2, a_2) + \dots}_{\Delta} \right] \right]$$

Bellman Equation:

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

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

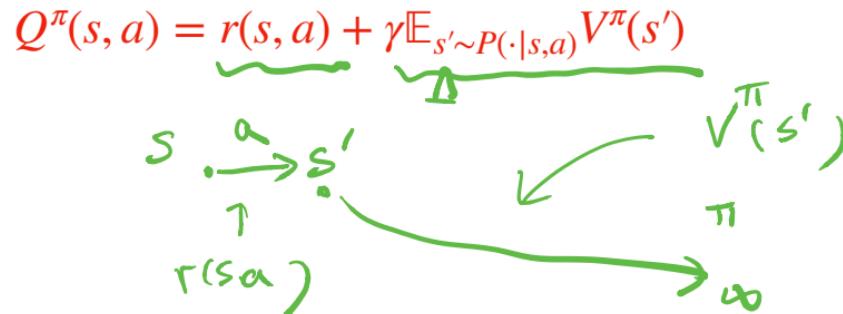
$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \middle| (\underbrace{s_0, a_0}_{(s, a)}, \underbrace{a_h}_{\pi(s_h)}, s_{h+1}) \sim P(\cdot | s_h, a_h) \right]$$

Bellman Equation:

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

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

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



Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^* : S \mapsto A, \text{ s.t., } V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

~~induction~~ induction \checkmark $\sqrt{\pi^*}$ $\pi^*(s) = a$

Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^* : S \mapsto A, \text{ s.t., } V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

We denote $\underbrace{V^*}_{\sim} := V^{\pi^*}, \underbrace{Q^*}_{\sim} := Q^{\pi^*}$

Optimal Policy

For infinite horizon discounted MDP, there exists a deterministic stationary policy

$$\pi^* : S \mapsto A, \text{ s.t., } V^{\pi^*}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.4 in the RL monograph]

We denote $V^* := V^{\pi^*}, Q^* := Q^{\pi^*}$

Theorem 1: Bellman Optimality

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$\Leftrightarrow r(s, a) + \gamma \mathbb{E}_{s' \sim P(s|a)} V^*(s')$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^\star(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^\star(s), \forall s$

$$\underbrace{V^\star(s)}_{\text{Bellman Eqn}} = r(s, \pi^\star(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^\star(s))} V^\star(s')$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^\star(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^\star(s), \forall s$

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

← Bell-Eqn

$= \underline{Q^\star(s, a)}$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^\star(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^\star(s), \forall s$

$$\begin{aligned} V^\star(s) &= r(s, \pi^\star(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^\star(s))} V^\star(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^\star(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \underline{V^\star(s')} \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[\underbrace{r(s', \pi^\star(s'))}_{\Delta} + \gamma \mathbb{E}_{s'' \sim P(s', \pi^\star(s'))} V^\star(s'') \right] \\ &\quad \text{max } \underset{a}{\Delta} (\Rightarrow \hat{\pi}(s')) \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^\star(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^\star(s), \forall s$

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^\star(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^\star(s), \forall s$

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s')$$

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

$$= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right]$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right]$$

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

$$\leq \mathbb{E} [r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots] = V^{\hat{\pi}}(s) \quad \hat{V}(s) \leq V^*(s) \leq \hat{V}(s) \Leftrightarrow \hat{V}(s) = V^*(s), \forall s$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^\star(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^\star(s), \forall s$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^*(s), \forall s$

This implies that $\arg \max_a Q^*(s, a)$ is an optimal policy

$$\hat{\pi} \quad V^{\hat{\pi}}(s) = V^*(s), \forall s \in S$$

Proof of Bellman Optimality

Theorem 2:

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

Proof of Bellman Optimality

Theorem 2:

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

$\forall s \in S$

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

Proof of Bellman Optimality

Theorem 2:

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

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

Proof of 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 $V(s) = V^*(s), \forall s$

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

Proof of Bellman Optimality

Theorem 2:

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

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

Recursivity

Proof of Bellman Optimality

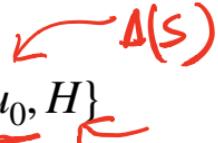
Theorem 2:

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

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s \sim P(s, a)} V^*(s')) \right| \\ &\stackrel{\text{O}}{=} \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} |V(s'') - V^*(s'')| \right) \\ &\leq \max_{a_1, a_2, \dots, a_{k-1}} \gamma^k \mathbb{E}_{s_k} |V(s_k) - V^*(s_k)| \quad k \rightarrow \infty \quad \gamma^k \rightarrow 0 \end{aligned}$$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \underline{\mu_0}, H\}$$



$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r(s_0, a_0), s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), r(s_{H-1}, a_{H-1})$$

$a_0 \dashrightarrow a_{H-1}$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), \underbrace{r(s_0, a_0)}_{\text{red arrow}}, s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), \underbrace{r(s_{H-1}, a_{H-1})}_{\text{red arrow}}$$

Objective function: $J(\pi) = \mathbb{E} \left[\sum_{h=0}^{H-1} r(s_h, a_h) \right]$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r(s_0, a_0), s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), r(s_{H-1}, a_{H-1})$$

Time-dependent value/Q function:

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{t=h}^{H-1} r(s_t, a_t) \mid s_h = s, a_t \sim \pi(s_t) \right], \quad Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{t=h}^{H-1} r(s_t, a_t) \mid s_h = s, a_h = a, a_t \sim \pi(s_t) \right]$$



Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r(s_0, a_0), s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), r(s_{H-1}, a_{H-1})$$

Time-dependent value/Q function:

$$V_h^\pi(s) = \mathbb{E} \left[\sum_{t=h}^{H-1} r(s_t, a_t) \mid s_h = s, a_t \sim \pi(s_t) \right], \quad Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{t=h}^{H-1} r(s_t, a_t) \mid s_h = s, a_h = a, a_t \sim \pi(s_t) \right]$$

$$V_h^\pi(s_h) = \mathbb{E}_{a_h \sim \pi(s_h)} \left[r(s_h, a_h) + \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} V_{h+1}^\pi(s_{h+1}) \right]$$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r(s_0, a_0), s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), r(s_{H-1}, a_{H-1})$$

Time-dependent optimal policy $\pi^\star = \{\pi_0^\star, \dots, \pi_{H-1}^\star\}$:

$$\pi_h^* : S \rightarrow \Delta(A)$$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r(s_0, a_0), s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), r(s_{H-1}, a_{H-1})$$



Time-dependent optimal policy $\pi^\star = \{\pi_0^\star, \dots, \pi_{H-1}^\star\}$:

$$Q_{H-1}^\star(s, a) = \underbrace{r(s, a)}_{=r(s,a)}, \pi_{H-1}^\star(s) = \arg \max_a Q_{H-1}^\star(s, a), V_{H-1}^\star(s) = \max_a Q_{H-1}^\star(s, a)$$

Finite Horizon Setting

$$\mathcal{M} = \{S, A, P, r, \mu_0, H\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad H \in \mathbb{N}^+, \quad \mu_0 \in \Delta(S)$$

Given $\pi : S \mapsto \Delta(A)$

$$s_0 \sim \mu_0, a_0 \sim \pi(s_0), r(s_0, a_0), s_1 \sim P(\cdot | s_0, a_0), \dots, s_{H-1} \sim P(\cdot | s_{H-2}, a_{H-2}), a_{H-1} \sim \pi(s_{H-1}), r(s_{H-1}, a_{H-1})$$

Time-dependent optimal policy $\pi^\star = \{\pi_0^\star, \dots, \pi_{H-1}^\star\}$:

$$Q_{H-1}^\star(s, a) = r(s, a), \pi_{H-1}^\star(s) = \arg \max_a Q_{H-1}^\star(s, a), V_{H-1}^\star(s) = \max_a Q_{H-1}^\star(s, a)$$

$$s \xrightarrow{h} a \xrightarrow{h+1} \underbrace{\pi^\star}_{Q_h^\star(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} V_{h+1}^\star(s')}$$

State (action) Occupancy

$\mathbb{P}_h(s; s_0, \pi)$: probability of π visiting s at time step $h \in \mathbb{N}$, starting at s_0

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

$$V^\pi(s_0) = \frac{1}{1 - \gamma} \sum_{s,a} d_{s_0}^\pi(s) \pi(a | s) r(s, a)$$

Valid Distribution

$s_0, a_0, s_1, a_1, \dots, s_n$

$$\Pr(s_0, a_0, s_1, a_1, \dots, s_n) = \underbrace{\pi(a_0 | s_0)}_{\Pr(s_0 | s_0, a_0)} \underbrace{P(s_1 | s_0, a_0)}_{\Pr(s_1 | s_1, a_1)} \cdots \underbrace{\pi(a_{n-1} | s_{n-1})}_{\Pr(s_n | s_{n-1}, a_{n-1})}$$

$$\mathbb{P}_h(s; s_0, \pi) = \sum_{a_0, s_1, \dots, a_{n-1}, s_{n-1}} \Pr(s_0, a_0, \dots, s_h = s)$$

State (action) Occupancy

$\mathbb{P}_h(\underline{s}, \underline{a}; s_0, \pi)$: probability of π visiting $(\underline{s}, \underline{a})$ at time step $h \in \mathbb{N}$, starting at s_0

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

$$\underline{V}^\pi(s_0) = \frac{1}{1 - \gamma} \sum_{s,a} d_{s_0}^\pi(s, a) r(s, a)$$

$$= \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, \pi \right]$$

π^* , Given P .
Δ
Linear programming