## **Value Iteration**

## **Lucas Janson and Sham Kakade**

#### CS/Stat 184: Introduction to Reinforcement Learning Fall 2022

## Today

- Recap
- Today:
  - An Iterative Algorithm: Value Iteration
  - Visitation distributions

# Recap

#### **Infinite horizon Discounted Setting**

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$
$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$
$$\mathsf{Policy} \ \pi : S \mapsto A$$

Quantities that allow us to reason policy's long-term effect:

Value function 
$$V^{\pi}(s) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h}) \middle| s_{0} = s, \pi\right]$$
  
Q function  $Q^{\pi}(s, a) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h}) \middle| (s_{0}, a_{0}) = (s, a), \pi\right]$ 

#### **Bellman Consistency Equations:**

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

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

#### Summary so far:

Every discounted MDP has some deterministic optimal policy  $\pi^{\star}$ , that dominates all other policies, everywhere

 $V^{\pi^{\star}}(s) \ge V^{\pi}(s), \forall \pi, \forall s$ 

#### Summary so far:

Every discounted MDP has some deterministic optimal policy  $\pi^*$ , that dominates all other policies, everywhere

 $V^{\pi^{\star}}(s) \ge V^{\pi}(s), \forall \pi, \forall s$ 

So we have,  $V^{\star} = V^{\pi^{\star}}$  and  $Q^{\star} = Q^{\pi^{\star}}$ .

•  $V^{\star}$  satisfies Bellman Equations:  $V^{\star}(s) = \max_{a} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right]$ 



• 
$$V^*$$
 satisfies Bellman Equations:  
 $V^*(s) = \max_{a} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right]$   
• If V satisfies the Bellman Equations,  
 $V(s) = \max_{a} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s') \right]$ ,  
then  $V = V^*$ .

•The optimal policy is:

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

## Today: Value Iteration

## Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ , how can we (approximately) find  $\pi^*$ ?

#### **Example of Optimal Policy** $\pi^{\star}$

Consider the following **deterministic** MDP w/ 3 states & 2 actions



#### What about this one...



Let's design an algorithm that computes  $V^{\star}/Q^{\star}$  for any given MDP

## Can we efficiently compute in a optimal policy? (polynomial in |S|, |A|, and other relevant quantities) how big r > P? (by forg fores if r > O(STA))

specify P(55)5g)

- Suppose we can efficiently compute  $V^{\pi}(s)$  for any given  $\pi: S \mapsto A$ .
  - Brute force search would to find  $\pi^*$  would still take  $|A|^{|S|}$  time.
- Can we construct an interactive algorithm based on the BEs?
  Will it converge?
  - •What is the computation time to get an approximate solution?

#### **Detour: fix-point solution**

Consider  $x^* = f(x^*), \quad f: [a, b] \mapsto [a, b]$ 

A naive approach to find  $x^*$ : Initialize  $x^0 \in [a, b]$ , repeat:  $x^{t+1} = f(x^t)$ 

#### **Detour: fix-point solution**

A naive approach to find  $x^*$ : Initialize  $x^0 \in [a, b]$ , repeat:  $x^{t+1} = f(x^t)$ 

= f (41) (x)

If f is a contraction mapping, i.e.,  $\forall x, x', |f(x) - f(x')| \le \gamma |x - x'|$ , for some  $\gamma \in [0,1)$ , then:  $x^t \to x^*$ , as  $t \to \infty$  $\langle x^t - x^* \rangle = (f(x^{t-1}) - f(x^*)) | \le \gamma |x - x^*|$ 

 $V(s) \in \mathbb{R}$ Define Bellman Operator  $\mathcal{T}$ : Bellman Equations:  $V(s) = \max_{a} \left[ r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s') \right] \quad \forall r(s, a) = \int_{V(s, a)} V(s') \left[ \frac{1}{V(s)} \right]$ •Any function  $V: S \mapsto \mathbb{R}$  can also be viewed as a vector in  $V \in \mathbb{R}^{|S|}$ . • Define  $\mathcal{T}: \mathbb{R}^{|S|} \mapsto \mathbb{R}^{|S|}$ . where  $(\mathcal{T}V)(s) := \max_{a} \left[ r(s,a) + \mathcal{E}_{s' \sim P(s,a)} V(s') \right]$ • Bellman equations in terms of  $\mathcal{T}$ : 2-Bellman Operator  $\mathcal{T}V = V$ 



## Value Iteration Algorithm:

1. Initialization: 
$$V^0 : ||V^0||_{\infty} \in \left[0, \frac{1}{1-\gamma}\right]$$
  
2. Iterate until convergence:  $V^{t+1} \leftarrow \mathcal{T}V^t$ 

#### Guarantee of VI:

We will see this fix-point iteration converges, i.e.,  $V^t \to V^*$ , as  $t \to \infty$ 

# Alternative Version: Bellman Operator $\mathcal{T}$ on Q (HW2 Q2 is the Q-version of the Bellman Equations)

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

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

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

## What is the Per-Iteration Computational Complexity? TV(S) = max 2r(Sa) + YE V(S) } a SinPlsa)

- Making the update  $V^{t+1} \leftarrow \mathcal{T}V^t$  explicit:
  - q eR IST. 141 • Define  $Q^{t+1}$ :  $\forall s, a \ Q^{t+1}(s, a) = r(s, a) + \gamma \sum P(s' \mid s, a) V^t(s')$ *s′*∈*S*  $E_{s' \times P | s \in V} V^{t} (s^{-})$ •Set  $V^{t+1}$ :  $\forall s \ V^{t+1}(s) = \max Q^{t+1}(s, a)$

#### What is the Per-Iteration Computational Complexity?

- Making the update  $V^{t+1} \leftarrow \mathcal{T}V^t$  explicit:
- O(IS)) opentions • Define  $Q^{t+1}$ :  $\forall s, a \ Q^{t+1}(s, a) = r(s, a) + \gamma \sum P(s' \mid s, a) V^t(s')$  $s' \in S$ •Set  $V^{t+1}$ :  $\forall s \ V^{t+1}(s) = \max Q^{t+1}(s, a) \qquad \bigcirc (A) ) \circ \rho = - f_{ons}$

•What is the order of the number of basic arithmetic operations?

# What is the Per-Iteration Computational Complexity? $4 \sqrt{2}$

- Making the update  $V^{t+1} \leftarrow \mathcal{T}V^t$  explicit:
  - Define  $Q^{t+1}$ :  $\forall s, a \ Q^{t+1}(s, a) = r(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^t(s')$ • Set  $V^{t+1}$ :  $\forall s \ V^{t+1}(s) = \max_{a} Q^{t+1}(s, a)$  $C\left( \left[ \leq \right] \left[ A \right] \right)$

•What is the order of the number of basic arithmetic operations?  $O(|S|^2 |A|)$ 

#### With matrix multiplication?

- Making the update  $V^{t+1} \leftarrow \mathcal{T}V^t$  explicit:
  - Define  $Q^{t+1}$ :  $\forall s, a \ Q^{t+1}(s, a) = r(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^t(s')$ • Set  $V^{t+1}$ :  $\forall s \ V^{t+1}(s) = \max_{a} Q^{t+1}(s, a)$

#### With matrix multiplication?

- Making the update  $V^{t+1} \leftarrow \mathcal{T}V^t$  explicit:
- Define  $Q^{t+1}$ :  $\forall s, a \ Q^{t+1}(s, a) = r(s, a) + \gamma \sum P(s' | s, a) V^{t}(s')$  $P_{(5,n)}, s' = P(s'(5, q))$  $s' \in S$ •Set  $V^{t+1}$ .  $\forall s \ V^{t+1}(s) = \max Q^{t+1}(s, a)$ pours column. •In terms of matrix multiplication, let us view  $\vec{r}$  as a vector,  $r \in \mathbb{R}^{|S| \cdot |A|}$  and P as a matrix  $P \in \mathbb{R}^{|S| \cdot |A| \times |S|}$   $\searrow \quad \bigcirc \quad \overleftarrow{c} \quad \overleftarrow{\tau} \quad$

#### With matrix multiplication?

- Making the update  $V^{t+1} \leftarrow \mathcal{T}V^t$  explicit:
- Define  $Q^{t+1}$ :  $\gg \forall s, a \ Q^{t+1}(s, a) = r(s, a) + \gamma \sum P(s' \mid s, a) V^t(s')$  $s' \in S$ •Set  $V^{t+1}$ .  $\forall s \ V^{t+1}(s) = \max Q^{t+1}(s, a)$ •In terms of matrix multiplication, let us view r as a vector,  $r \in \mathbb{R}^{|S| \cdot |A|}$  and P as a matrix  $P \in \mathbb{R}^{|S| \cdot |A| \times |S|}$  $Q^{t+1} = r + \gamma P V^t$

## Outline:

1: An Iterative Algorithm: Value Iteration (a fix-point iteration algorithm again)

2: Convergence? How fast? (Via the <u>contraction</u> argument again!)

## Convergence of Value Iteration:

#### *Lemma [contraction]*: Given any V, V', we have: $\|\mathscr{T}V - \mathscr{T}V'\|_{\infty} \leq \gamma \|V - V'\|_{\infty}$

 $|(\mathcal{T}V)(s) - (\mathcal{T}V')(s)| = \left| \max_{a} \left\{ r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V(s') \right\} - \max_{a} \left\{ r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V'(s') \right\} \right|$ Proof: Max In(Salt & Francis V(sr) - (n(Salt& F V(Sr))) = & max [ E V(s) - V(s) - V(s) ]  $\leq \gamma \max_{a} E_{w} \left[ V(S^{r}) - V(S^{r}) \right]$ 

Convergence of Value Iteration: *Lemma [contraction]*: Given any Q, Q', we have:  $\|\mathscr{T}Q - \mathscr{T}Q'\|_{\infty} \leq \gamma \|Q - Q'\|_{\infty}$ 

**Proof:** 

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

## Convergence of Value Iteration:

*Lemma [Convergence]*: Given  $V^0$ , we have:  $\|V^t - V^{\star}\|_{\infty} \le \gamma^t \|V^0 - V^{\star}\|_{\infty}$ 

Proof:

$$\|V^{t} - V^{\star}\|_{\infty} = \|\mathscr{T}V^{t-1} - \mathscr{T}V^{\star}\|_{\infty} \le \gamma \|V^{t-1} - V^{\star}\|_{\infty}$$

## Convergence of Value Iteration:

*Lemma [Convergence]*: Given  $V^0$ , we have:  $\|V^t - V^\star\|_{\infty} \le \gamma^t \|V^0 - V^\star\|_{\infty}$ 

$$\|V^{t} - V^{\star}\|_{\infty} = \|\mathscr{T}V^{t-1} - \mathscr{T}V^{\star}\|_{\infty} \le \gamma \|V^{t-1} - V^{\star}\|_{\infty}$$

$$\ldots \leq \gamma^{t+\gamma} \|V^0 - V^\star\|_{\infty}$$

$$\begin{array}{c} \text{Computational Complexity of VI} \\ \begin{array}{c} & & \\$$

· SEE fixed The Quality of Policy: **Theorem:** For any V, let  $\pi^{\bullet}(s) = \arg \max_{a} \left[ r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s') \right]$ , then if  $v \in \mathcal{O}$  is  $V^{\pi^{\circ}}(s) \ge V^{\star}(s) - \frac{2\gamma}{1-\gamma} \|V - V^{\star}\|_{\infty} + \frac{MDP}{VP} + \frac{2\gamma}{1-\gamma} \|V - V^{\star}\|_{\infty}$  $E \mathcal{E} (FY)$ finds, a policy TT S.t. VI  $|V^{T} - V^{*}|_{to} \leq \mathcal{Z}$ in # iterstices log ( scirb)

## The Quality of Policy:

**Theorem:** For any V, let  $\pi^{\bullet}(s) = \arg \max_{a} \left[ r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s') \right]$ , then  $V^{\pi^{\bullet}}(s) \ge V^{\star}(s) - \frac{2\gamma}{1-\gamma} \|V - V^{\star}\|_{\infty}$ 

> see slides to, proof

(not necessary)

#### Proof:

• To prove the theorem, it suffice to show that:

$$\|V^{\pi^{\flat}} - V^{\star}\| \leq \frac{2\gamma}{1 - \gamma} \|V - V^{\star}\|_{\infty}$$

## Understanding the sampling

"Occupancy measures" are a helpful concept

Assume we start at  $s_0$ , following  $\pi$  to step h, what's probability of seeing a trajectory:

 $(s_0, a_0, s_1, a_1, \dots, s_h, a_h)$ ?

Assume we start at  $s_0$ , following  $\pi$  to step h, what's probability of seeing a trajectory:  $(s_0, a_0, s_1, a_1, \dots, s_h, a_h)$ ?

Let's write 
$$\pi$$
 as a delta distribution, i.e.,  $\pi(a \mid s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$ 

Assume we start at  $s_0$ , following  $\pi$  to step h, what's probability of seeing a trajectory:  $(s_0, a_0, s_1, a_1, \dots, s_h, a_h)$ ?

Let's write 
$$\pi$$
 as a delta distribution, i.e.,  $\pi(a \mid s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$ 



Assume we start at  $s_0$ , following  $\pi$  to step h, what's probability of seeing a trajectory:  $(s_0, a_0, s_1, a_1, \dots, s_h, a_h)$ ?

Let's write 
$$\pi$$
 as a delta distribution, i.e.,  $\pi(a \mid s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$ 



Assume we start at  $s_0$ , following  $\pi$  to step h, what's probability of seeing a trajectory:  $(s_0, a_0, s_1, a_1, \dots, s_h, a_h)$ ?

Let's write 
$$\pi$$
 as a delta distribution, i.e.,  $\pi(a \mid s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$ 



$$\mathbb{P}^{\pi}(s_0, a_0, \dots, s_h, a_h)$$

... =  $\pi(a_0 | s_0) P(s_1 | s_0, a_0) \pi(a_1 | s_1) P(s_2 | s_1, a_1) \dots P(s_h | s_{h-1}, a_{h-1}) \pi(a_h | s_h)$ 



Q: what's the probability of  $\pi$  visiting state (*s*,a) at time step h?

#### State-action distribution at time step h



$$\mathbb{P}^{\pi}(a_0,\ldots,s_h,a_h\,|\,s_0,\pi)$$

 $= \pi(a_0 | s_0) P(s_1 | s_0, a_0) \pi(a_1 | s_1) P(s_2 | s_1, a_1) \dots P(s_h | s_{h-1}, a_{h-1}) \pi(a_h | s_h)$ 

Q: what's the probability of  $\pi$  visiting state (*s*,a) at time step h?

$$\mathbb{P}_{h}^{\pi}(s_{h}, a_{h} \mid s_{0}, \pi) = \sum_{a_{0}, s_{1}, a_{1}, \dots, s_{h-1}, a_{h-1}} \mathbb{P}^{\pi}(a_{0}, \dots, s_{h-1}, a_{h-1} s_{h} = s, a_{h} = a \mid s_{0}, \pi)$$

Probability of  $\pi$  visiting (s, a) at h, starting from  $s_0$ 

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

Can you show that this is a valid distribution?

Probability of  $\pi$  visiting (s, a) at h, starting from  $s_0$ 

$$d_{s_0}^{\pi}(s,a) = (1-\gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s,a;s_0)$$
  
Can you show that this is a valid distribution?  
$$V^{\pi}(s_0) = \frac{1}{1-\gamma} \sum_{s,a} d_{s_0}^{\pi}(s,a) r(s,a)$$

s,a

Probability of  $\pi$  visiting (s, a) at h, starting from  $s_0$ 

$$d_{s_0}^{\pi}(s,a) = (1-\gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s,a;s_0)$$
  
Can you show that this is a valid distribution?  
$$V^{\pi}(s_0) = \frac{1}{1-\gamma} \sum_{s,a} d_{s_0}^{\pi}(s,a) r(s,a)$$

Can you show the above is true?

Probability of  $\pi$  visiting (s, a) at h, starting from  $s_0$ 

$$d_{s_0}^{\pi}(s,a) = (1-\gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s,a;s_0)$$
  
Can you show that this is a valid distribution?  
$$V^{\pi}(s_0) = \frac{1}{1-\gamma} \sum_{s,a} d_{s_0}^{\pi}(s,a)r(s,a)$$
 HW0 questions

Can you show the above is true?

Probability of  $\pi$  visiting (s, a) at h, starting from  $s_0$ 

$$\mathbb{P}_{h}^{\pi}(s_{h}, a_{h} | s_{0}, \pi) = \sum_{a_{0}, s_{1}, a_{1}, \dots, s_{h-1}, a_{h-1}} \mathbb{P}^{\pi}(s_{0}, a_{0}, \dots, s_{h-1}, a_{h-1}s_{h} = s, a_{h} = a)$$

$$d_{s_{0}}^{\pi}(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^{h} \mathbb{P}_{h}^{\pi}(s, a; s_{0})$$
Can you show that this is a valid distribution?

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

**HW0** questions!

Can you show the above is true?

#### **Summary Today**

- Value iteration: an iterative algorithm with a "linear" convergence rate.
- The concept of an "occupancy measure".

1-minute feedback form: <a href="https://bit.ly/3RHtlxy">https://bit.ly/3RHtlxy</a>

