# Computation Limits and The LP Formulation

H: )

## Sham Kakade and Wen Sun

**CS 6789: Foundations of Reinforcement Learning** 

## Today:

- Next two lectures:
  - Today: computational complexity & the linear programming approach
  - Next lecture: statistical complexity (when the MDP is not known)

## Today:

- Next two lectures:
  - Today: computational complexity & the linear programming approach
  - Next lecture: statistical complexity (when the MDP is not known)
- Today: Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$  can we exactly compute  $Q^*$  (or find  $\pi^*$ ) in polynomial time?

## Today:

- Next two lectures:
  - Today: computational complexity & the linear programming approach
  - Next lecture: statistical complexity (when the MDP is not known)
- Today: Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$  can we exactly compute  $Q^{\star}$  (or find  $\pi^{\star}$ ) in polynomial time?
- But first, our recap:
  - value/policy iteration + contraction

## Recap

#### Define Bellman Operator $\mathcal{T}$ :

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

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

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

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

#### Policy Iteration Algorithm:

Closed-form for PE (see AJKS) 2. Policy Evaluation:  $\pi^0 : S \mapsto \Delta(A)$ 2. Policy Evaluation:  $Q^{\pi^t}(s, a), \forall s, a$ 3. Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$ 

## Final Quality of the Policy (for VI):

• 
$$\pi^{t}$$
:  $\pi^{t}(s) = \arg \max_{a} Q^{t}(s, a)$   
Theorem:  $V^{\pi^{t}}(s) \ge V^{\star}(s) - \frac{2\gamma^{t}}{1-\gamma} \|Q^{0} - Q^{\star}\|_{\infty} \forall s \in S$  •  $\bigvee \mathcal{I}$ 

#### Final Quality of the Policy (for VI):

• 
$$\pi^t : \pi^t(s) = \arg \max_a Q^t(s, a)$$
  
Theorem:  $V^{\pi^t}(s) \ge V^{\star}(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^{\star}\|_{\infty} \forall s \in S$ 

• Corollary: Set 
$$Q^0 = 0$$
. After  $t \ge \frac{\log \frac{2}{\epsilon(1-\gamma)^2}}{1-\gamma}$  iterations, we have:  
 $V^{\pi^t}(s) \ge V^{\star}(s) - \epsilon \quad \forall s \in S$ 

#### Final Quality of the Policy (for VI):

• 
$$\pi^t : \pi^t(s) = \arg \max_a Q^t(s, a)$$
  
Theorem:  $V^{\pi^t}(s) \ge V^{\star}(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^{\star}\|_{\infty} \forall s \in S$ 

• Corollary: Set 
$$Q^0 = 0$$
. After  $t \ge \frac{\log \frac{2}{\epsilon(1-\gamma)^2}}{1-\gamma}$  iterations, we have:  
 $V^{\pi^t}(s) \ge V^{\star}(s) - \epsilon \quad \forall s \in S$ 

• Same rate for PI.

# Today

- Complexity to compute an **exact solution** given *M*. (Aside: why?)
  - Assume basic arithmetic operations  $(+, -, x, \div)$  take unit time.

- Complexity to compute an **exact solution** given *M*. (Aside: why?)
  - Assume basic arithmetic operations  $(+,-,x,\div)$  take unit time.
- Polytime computation: Suppose that  $(P, r, \gamma)$  in our MDP  $\mathcal{M}$  is specified with rational entries, where  $L(P, r, \gamma)$  is total bit-size required to specify  $(P, r, \gamma)$ .
  - Can we (exactly) compute  $Q^*$  in time poly  $(S, A, L(P, r, \gamma))$ ?

- Complexity to compute an **exact solution** given *M*. (Aside: why?)
  - Assume basic arithmetic operations  $(+,-,x,\div)$  take unit time.
- Polytime computation: Suppose that  $(P, r, \gamma)$  in our MDP  $\mathcal{M}$  is specified with rational entries, where  $L(P, r, \gamma)$  is total bit-size required to specify  $(P, r, \gamma)$ .
  - Can we (exactly) compute  $Q^*$  in time poly  $(S, A, L(P, r, \gamma))$ ?
- (nearly) strongly polynomial time: Suppose (*P*, *r*) is specified with real numbers. Can we compute *Q*<sup>\*</sup> in poly(*S*, *A*, log(1/(1 - γ)), with no dependence on *L*(*P*, *r*, γ)?

- Complexity to compute an **exact solution** given *M*. (Aside: why?)
  - Assume basic arithmetic operations  $(+, -, x, \div)$  take unit time.
- Polytime computation: Suppose that  $(P, r, \gamma)$  in our MDP  $\mathcal{M}$  is specified with rational entries, where  $L(P, r, \gamma)$  is total bit-size required to specify  $(P, r, \gamma)$ .
  - Can we (exactly) compute  $Q^{\star}$  in time poly  $(S, A, L(P, r, \gamma))$ ?
- (nearly) strongly polynomial time: Suppose (P, r) is specified with real numbers. Can we compute  $Q^{\star}$  in poly  $(S, A, \log(1/(1 - \gamma)))$ , with no dependence on  $L(P, r, \gamma)$ ? with S, A, +

• Scalings:

• How does the complexity scale with the "horizon"  $1/(1 - \gamma)$ ? With  $L(P, r, \gamma)$ ?

Computational Complexities of our Iterative Algorithms

When the sub-optimality gap between Q<sup>t</sup> and Q<sup>\*</sup> is less than 2<sup>-L(P,r,γ)</sup>, than the greedy policy will be optimal.
 (by a standard argument in optimization)

- When the sub-optimality gap between Q<sup>t</sup> and Q<sup>\*</sup> is less than 2<sup>-L(P,r,γ)</sup>, than the greedy policy will be optimal.
   (by a standard argument in optimization)
- VI:
  - $O(\log(1/(\epsilon(1-\gamma))/(1-\gamma))$  iterations suffice for an  $\epsilon$ -accurate solution.
  - Per iteration complexity:  $S^2A$

- When the sub-optimality gap between Q<sup>t</sup> and Q<sup>\*</sup> is less than 2<sup>-L(P,r,γ)</sup>, than the greedy policy will be optimal.
   (by a standard argument in optimization)
- VI:
  - $O(\log(1/(\epsilon(1-\gamma))/(1-\gamma)))$  iterations suffice for an  $\epsilon$ -accurate solution.
  - Per iteration complexity:  $S^2A$
- Poly runtime? For fixed  $\gamma$ , YES • VI is poly with  $O\left(S^2A \cdot \frac{L(P, r, \gamma) \log(1/(1 - \gamma))}{1 - \gamma}\right)$  complexity.

- When the sub-optimality gap between Q<sup>t</sup> and Q<sup>\*</sup> is less than 2<sup>-L(P,r,γ)</sup>, than the greedy policy will be optimal.
   (by a standard argument in optimization)
- VI:
  - $O(\log(1/(\epsilon(1-\gamma))/(1-\gamma)))$  iterations suffice for an  $\epsilon$ -accurate solution.
  - Per iteration complexity:  $S^2A$
- Poly runtime? For fixed  $\gamma$ , YES • VI is poly with  $O\left(S^2A \cdot \frac{L(P, r, \gamma) \log(1/(1 - \gamma))}{1 - \gamma}\right)$  complexity.
- Strongly poly? NO (There are counterexamples)

## **Policy Iteration**

VI and checked

- PI Per iteration complexity:  $S^3 + S^2 A$  (why?) complex  $Q^{\pi}$ 
  - PI is more costly than VI per iteration.
  - PI is observed to be much faster than VI to obtain an exact opt policy.

## **Policy Iteration**

- PI Per iteration complexity:  $S^3 + S^2A$  (why?)
  - PI is more costly than VI per iteration.
  - PI is observed to be much faster than VI to obtain an exact opt policy.
- Poly runtime? • For fixed  $\gamma$ , YES • Pl is  $O\left((S^3 + S^2A) \frac{L(P, r, \gamma) \log(1/(1 - \gamma))}{1 - \gamma}\right)$

## **Policy Iteration**

- PI Per iteration complexity:  $S^3 + S^2A$  (why?)
  - PI is more costly than VI per iteration.
  - PI is observed to be much faster than VI to obtain an exact opt policy.
- Poly runtime?
  - For fixed  $\gamma$ , YES

• Pl is 
$$O((S^3 + S^2 A) \frac{L(P, r, \gamma) \log(1/(1 - \gamma))}{1 - \gamma})$$

- Strongly poly?
  - Does PI compute an optimal policy in time independent of  $L(P, r, \gamma)$ ? (ignoring other dependencies?)

## Is PI a strongly poly algo?

- Does PI compute an optimal policy in time independent of  $L(P, r, \gamma)$ ?
  - Yes: after  $A^S$  iterations Why? There are at most  $A^S$  policies, and PI is monotonic.
  - Refinement: [Mansour & Singh '99] PI halts after  $A^S/S$  iterations.

## Is PI a strongly poly algo?

- Does PI compute an optimal policy in time independent of  $L(P, r, \gamma)$ ?
  - Yes: after  $A^S$  iterations Why? There are at most  $A^S$  policies, and PI is monotonic.
  - Refinement: [Mansour & Singh '99] PI halts after  $A^S/S$  iterations.
- Is PI strongly polynomial?

For fixed  $\gamma$ , yes: [Ye '12] PI halts after  $\frac{S^2 A \log(S^2/(1-\gamma))}{1-\gamma}$  iterations.

## Summary Table

	Value Iteration	Policy Iteration	LP-based Algorithms
Poly.	$S^2 A \frac{L(P,r,\gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$(S^3 + S^2 A) \frac{L(P,r,\gamma)\log\frac{1}{1-\gamma}}{1-\gamma}$	?
Strongly Poly.	Х	$\left(S^3 + S^2 A\right) \cdot \min\left\{\frac{A^S}{S}, \frac{S^2 A \log \frac{S^2}{1-\gamma}}{1-\gamma}\right\}$	?

- VI Per iteration complexity:  $S^2A$
- PI Per iteration complexity:  $S^3 + S^2A$

Are VI and PI Polynomial Time algorithms? (technically, no)

Is there a poly (and strongly poly) time algo for an MDP? YES! Linear Programming

• We can write the Bellman equations with values rather than Q-values:  $V(s) = \max_{a} \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V(s) \right] \right\}$ 

- We can write the Bellman equations with values rather than Q-values:  $V(s) = \max_{a} \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V(s) \right] \right\}$
- An equivalent way to write the Bellman equations is as a linear program.

- We can write the Bellman equations with values rather than Q-values:  $V(s) = \max_{a} \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V(s) \right] \right\} \qquad \begin{array}{c} \downarrow \varsigma_{a} & \downarrow & \downarrow \\ \varsigma_{c} \rightarrow & \downarrow \\ \varsigma_{c} \rightarrow$
- An equivalent way to write the Bellman equations is as a linear program.
- With variables  $V \in \mathbb{R}^{S}$ , the LP is:  $\min V(s_{0})$   $V \in \mathbb{R}^{S}$ , the LP is:  $\min \frac{1}{5} \stackrel{\frown}{\downarrow} \stackrel{\downarrow}{\downarrow} \stackrel{\downarrow}{\downarrow} \stackrel{\frown}{\downarrow} \stackrel{\frown}{\downarrow} \stackrel{\frown}{\downarrow} \stackrel{\downarrow}{\downarrow} \stackrel{\downarrow}$

s.t.  $V(s) \ge r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} V(s') \quad \forall s, a \in S \times A$ 

#### LP Runtimes and Comments

• An LP solver gives us a poly time algorithm.

## LP Runtimes and Comments

- An LP solver gives us a poly time algorithm.
- [Ye, '05]: there is an interior point algorithm (CIPA) which is ("nearly") strongly polynomial.

## LP Runtimes and Comments

- An LP solver gives us a poly time algorithm.
- [Ye, '05]: there is an interior point algorithm (CIPA) which is ("nearly") strongly polynomial.
- Comments:
  - VI is best thought of as a fixed point algorithm
  - PI is equivalent to a (block) simplex algorithm (Recall the simplex algo, in general, could be exp time. But not for MDPS, at least for fixed *γ*.)

## Summary Table

	Value Iteration	Policy Iteration	LP-based Algorithms
Poly.	$S^2 A \frac{L(P,r,\gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$(S^3 + S^2 A) \frac{L(P,r,\gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$S^3AL(P,r,\gamma)$
Strongly Poly.	Х	$\left(S^3 + S^2 A\right) \cdot \min\left\{\frac{A^S}{S}, \frac{S^2 A \log \frac{S^2}{1-\gamma}}{1-\gamma}\right\}$	$S^4 A^4 \log \frac{S}{1-\gamma}$

CTAPA

- VI Per iteration complexity:  $S^2A$
- PI Per iteration complexity:  $S^3 + S^2A$
- The LP approach is only logarithmic in  $1 \gamma$
- The linear programming is helpful in understanding the problem. (even though it is not used often)

#### What about the Dual LP?

#### What about the Dual LP?

- Let us now consider the dual LP.
  - It is also very helpful conceptually.
  - In some cases, it also provides a reasonable algorithmic approach

#### What about the Dual LP?

- Let us now consider the dual LP.
  - It is also very helpful conceptually.
  - In some cases, it also provides a reasonable algorithmic approach

• Let us start by understanding the dual variables and the "state-action polytope"

# State-Action Visitation Measures

• For a fixed (possibly stochastic) policy  $\pi$ , define the Pr(So, ao rom state-action visitation distribution  $\nu^{\pi}$  as:  $\nu_{\infty}^{(\pi)}(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} \Pr^{\pi}(s_{t} = s, a_{t} = a \mid s_{0}) \quad \forall (a) \quad \forall s_{0}, \tau )$ where  $\Pr^{\pi}(s_t = s, a_t = a \mid s_0)$  is the state-action visitation probability when we execute  $\pi$  starting at state  $s_0$ . think it as rector 

#### **State-Action Visitation Measures**

• For a fixed (possibly stochastic) policy  $\pi$ , define the state-action visitation distribution  $\nu^{\pi}$  as:

$$\nu_{g}^{\pi}(s,a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^{t} \Pr^{\pi}(s_{t} = s, a_{t} = a \mid s_{0})$$

where  $\Pr^{\pi}(s_t = s, a_t = a \mid s_0)$  is the state-action visitation probability when we execute  $\pi$  starting at state  $s_0$ .

• We can verify that have  $\nu^{\pi}$  satisfies, for all states  $s \in S$ :  $\sum_{\substack{a \\ r}} \nu_{\underline{s}}^{\pi}(s, a) = (1 - \gamma)I(s = s_0) + \gamma \sum_{\substack{s', a'}} P(s \mid s', a') \nu_{\underline{s}}^{\pi}(s', a')$ 

### The "State-Action" Polytope

• Let us define the state-action polytope K as follows:  $K := \begin{cases} \nu \mid \nu \ge 0 \text{ and} \end{cases}$ 

$$\sum_{a} \nu(s, a) = (1 - \gamma)I(s = s_0) + \gamma \sum_{s', a'} P(s \mid s', a')\nu(s', a')$$

## The "State-Action" Polytope

• Let us define the state-action polytope K as follows:  $K := \begin{cases} \nu \mid \nu \ge 0 \text{ and} \end{cases}$ 

$$\sum_{a} \nu(s, a) = (1 - \gamma)I(s = s_0) + \gamma \sum_{s', a'} P(s \mid s', a')\nu(s', a') \bigg\}$$

• This set precisely characterizes all state-action visitation distributions:

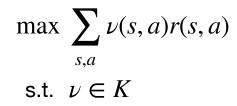
## The "State-Action" Polytope

• Let us define the state-action polytope K as follows:

$$K := \left\{ \nu \mid \nu \ge 0 \text{ and} \right.$$
$$\sum_{a} \nu(s, a) = (1 - \gamma)I(s = s_0) + \gamma \sum_{s', a'} P(s \mid s', a')\nu(s', a') \right\}$$

• This set precisely characterizes all state-action visitation distributions: Lemma:  $\nu \in K$  if and only if there exists a (possibly randomized) policy  $\pi$  s.t.  $\nu^{\pi} = \nu$ 

#### The Dual LP



- One can verify that this is the dual of the primal LP.
- Note that K is a polytope