Computational Limits & The LP formulation

Announcements

- HW0: due this Thursday 11:59pm
- Gradescope (please self-enroll)
- Norms for class?
 - Video:
 - Questions:

Today:

- Recap:
 - value/policy iteration + contraction
- Today: computational complexity & the linear programming approach

Question: Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ can we exactly compute Q^* (or find π^*) in polynomial time?

Recap

Define Bellman Operator \mathcal{T} :

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

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

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

Value Iteration Algorithm:

1. Initialization:
$$Q^0 : ||Q^0||_{\infty} \in (0, \frac{1}{1-\gamma})$$

2. Iterate until convergence: $Q^{t+1} = \mathcal{T}Q^t$

Policy Iteration Algorithm:

Closed-form for PE 1. Initialization: $\pi^0 : S \mapsto \Delta(A)$ (see 1.1.3 in Monograph) 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):

 π^t : $\pi^t(s) = \arg \max Q^t(s, a)$ Theorem: $V^{\pi^t}(s) \ge V^{\star}(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^{\star}\|_{\infty} \forall s \in S$ $\gamma^{t} = (1 - (1 - \gamma))^{t}$ $= (1 - \gamma) \cdot t$ $\leq e$ $V^{\pi'}(s) \ge V^{\star}(s) - \epsilon \quad \forall s \in S$ Same rate for Pl.

Today

Polynomial Time & Strongly Polynomial Complexity

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

Polynomial Time & Strongly Polynomial Complexity

- Complexity to compute an exact solution given \mathcal{M} . (Aside: Why?
 - Assume that basic arithmetic operations $(+,-,x,\div)$ take unit time.
- Polytime computation: Suppose that (P, r, γ) in our MDP \mathscr{M} is specified with rational entries, where $L(P, r, \gamma)$ is total bit-size required to specify (P, r, γ) . $\mathbb{P}\left(S = \# | S = 2, q = 8 \right) = \frac{10}{11} \quad (-) = 10^{-4}$ Can we (exactly) compute Q^* in time that is polynomial in $L(P, r, \gamma)$, # states *S*, and #actions *A*.

Polynomial Time & Strongly Polynomial Complexity

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

Can we (exactly) compute Q^* in time that is polynomial in $L(P, r, \gamma)$, # states *S*, and #actions *A*.

 Strongly polynomial time: Suppose (*P*, *r*, *γ*) is specified with real numbers. Can we compute *Q*[★] in poly(*S*, *A*, log(1/(1 − *γ*))), with no dependence on *L*(*P*, *r*, *γ*)? Computational Complexities of our Iterative Algorithms

 When the gap in the current objective value and the optimal objective value is smaller than 2^{-L(P,r,γ)}, then the greedy policy will be optimal. (this is a standard argument in optimization)

- When the gap in the current objective value and the optimal objective value is smaller than 2^{-L(P,r,γ)}, then the greedy policy will be optimal. (this is a standard argument in optimization)
 VI:
 - needs $\log(1/(\epsilon(1-\gamma))/(1-\gamma))$ iterations to obtain an ϵ accurate solution.
 - Per iteration complexity: S^2A

- When the gap in the current objective value and the optimal objective value is smaller than $2^{-L(P,r,\gamma)}$, then the greedy policy will be optimal. (this is a standard argument in optimization)
- VI:

• needs $\log(1/(\epsilon(1-\gamma))/(1-\gamma))$ iterations to obtain an ϵ accurate solution.

- Per iteration complexity: S^2A
- y: Technically, not poly. (For tixed & itis!) • Poly runtime? For fixed γ , VI is poly:

$$\int S^2 A \frac{L(P, r, \gamma) \log(1/(1-\gamma))}{1-\gamma}$$

- When the gap in the current objective value and the optimal objective value is smaller than 2^{-L(P,r,γ)}, then the greedy policy will be optimal. (this is a standard argument in optimization)
- VI:

• needs $\log(1/(\epsilon(1-\gamma))/(1-\gamma))$ iterations to obtain an ϵ accurate solution.

- Per iteration complexity: S^2A
- Poly runtime? For fixed γ , VI is poly:

$$S^{2}A \frac{L(P, r, \gamma) \log(1/(1-\gamma))}{1-\gamma}$$

Strongly poly? No must have dop. on $L(P, r, \beta)$

Compste, TI Policy Iteration to compte V tor jiren a involver solving 2. a linear system.

NT

- PI Per iteration complexity: $S^3 + S^2A$
 - 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: $S^3 + S^2A$
 - PI is more costly than VI per iteration.
 - PI is observed to be much faster than VI to obtain an exact opt policy.

(upper band)

• Poly runtime? For fixed γ , \forall is poly:

$$\underbrace{(S^3 + S^2 A)}_{\text{A} \text{ formula}} \underbrace{\frac{L(P, r, \gamma) \log(1/(1 - \gamma))}{1 - \gamma}}_{1 - \gamma}$$

Policy Iteration

- PI Per iteration complexity: $S^3 + S^2A$
 - 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 γ , VI is poly:

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

• Does PI compute an optimal policy in time independent of $L(P, r, \gamma)$? les l finite # of Policies

Is PI a strongly poly algo?

Does PI compute an optimal policy in time independent of *L*(*P*, *r*, *γ*)?
Yes: after *A^S* iterations (*A^S* is the number of policies) 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*, *γ*)?
 Yes: after *A^S* iterations (*A^S* is the number of policies) Refinement: [Mansour & Singh '99] PI halts after *A^S/S* iterations.
- Is PI strongly polynomial? For fixed γ , 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 | | | | | |
| these are | • PI Per ite | eration complexity: $S^3 + S^2 A$ | | | |

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

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

The Primal Linear Program

- We can write the Bellman equations with values rather than Q-values: $V(s) = \max_{a} \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[V(s) \right] \right\}$
- An equivalent way to write the Bellman equations is as a linear program.
 With variables V ∈ ℝ^S, the LP is:
 - $\min V(s_0)$ s.t. $V(s) \ge r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} V(s') \quad \forall s, a \in S \times A \quad \text{feas: ble}$ $\sum_{s'} P(s' \mid s, a) \vee (s') \quad \forall s, a \in S \times A \quad \text{feas: ble}$

LP Runtimes and Comments f(t,s)

• Using a polytime LP solver, gives us a poly time algorithm.

LP Runtimes and Comments

Using a polytime LP solver, gives us a poly time algorithm.
[Ye, '05]: there is an interior point algorithm (CIPA) which is also strongly polynomial.
4 (1/1-3) and 10 dep.
M L (4r,r)

LP Runtimes and Comments

- Using a polytime LP solver, gives us a poly time algorithm.
- [Ye, '05]: there is an interior point algorithm (CIPA) which is also strongly polynomial.
- Relations:
 - 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 | $\left(S^3 + S^2 A\right) \cdot \min\left\{\frac{A^S}{S}, \frac{S^2 A \log \frac{S^2}{1-\gamma}}{1-\gamma}\right\}$ | $S^4 A^4 \log \frac{S}{1-\gamma}$ |



- VI Per iteration complexity: S^2A
- PI Per iteration complexity: $S^3 + S^2A$
- The LP approach is only logarithmic in $1-\gamma$

What about the Dual LP?

- The linear programming is helpful in understanding the problem. (even though it is not used often)
- 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

State-Action Visitation Measures

• For a fixed (possibly stochastic) policy π , define the state-action visitation distribution ν^{π} as:

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

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

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

• Let us define the state-action polytope K as follows: $K := \left\{ \nu \mid \nu \ge 0 \text{ and} \right.$ $\sum_{a} \nu(s, a) = (1 - \gamma)I(s = s_0) + \gamma \sum_{s', a'} P(s \mid s', a')\nu(s', a') \right\}$

• Let us define the state-action polytope K as follows: $K := \left\{ \nu \mid \nu \ge 0 \text{ and} \right.$ $\sum_{a} \nu(s, a) = (1 - \gamma)I(s = s_0) + \gamma \sum_{s', a'} P(s \mid s', a')\nu(s', a') \right\}$

- 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:

- 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: K is equal to the set of all feasible state-action distributions, i.e. $\nu \in K$ if and only if there exists a (possibly randomized) policy π 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