Lecture 9 Planning with global access and uniform realizability: Part II

Global access

Can get all the features at all states, can preprocess it

$\underline{\zeta}$ –uniform action-value realizability

$$\sup_{\pi \in \text{DETML}} \inf_{\theta} \|q^{\pi} - \Phi\theta\|_{\infty} \leq \zeta$$

least-squares policy evaluation

Least-squares Policy Evaluation

1. Rollouts from a set C of wellchosen state-action pairs

- 2. Average over those
- 3. Least-squares fit



 $(\mathbb{E}\big[\hat{R}_m(z)\big] = q^\pi(z))$

$$\hat{ heta} = rg\min_{ heta \in \mathbb{R}^d} \sum_{z \in \mathcal{C}} arrho(z) \Big(\langle heta, arphi(z)
angle - \hat{R}_m(z) \Big)^2$$

$$\hat{R}_m(z) = rac{1}{m}\sum_{j=1}^m\sum_{t=0}^{H^{(j)}-1}r_{A_t^{(j)}}(S_t^{(j)})$$

 $\mathbb{P}(H^{(j)} \geq t+1) = \gamma^t$

First head with $X_1, X_2, ...$: $\mathbb{P}(X_i = \text{Head}) = 1 - \gamma$

$$\hat{ heta} = rg\min_{ heta \in \mathbb{R}^d} \sum_{z \in \mathcal{C}} arrho(z) \Big(\langle heta, arphi(z)
angle - \hat{R}_m(z) \Big)^2$$

Lemma (LSPE-*G* **extrapolation error control):** Fix any full-rank feature-map $\varphi : \mathbb{Z} \to \mathbb{R}^d$ and take the set $\mathcal{C} \subset \mathbb{Z}$ and the weighting function $\varrho \in \Delta_1(\mathcal{C})$ as in the Kiefer-Wolfowitz theorem. Fix an arbitrary policy π and let θ and ε_{π} such that $q^{\pi} = \Phi \theta + \varepsilon_{\pi}$ and assume that immediate rewards belong to the interval [0, 1]. Let $\hat{\theta}$ be as in Eq. (6). Then, for any $0 \leq \delta \leq 1$, with probability $1 - \delta$,

$$\left\|q^{\pi} - \Phi\hat{\theta}\right\|_{\infty} \leq \|\varepsilon_{\pi}\|_{\infty}(1 + \sqrt{d}) + \sqrt{d}\left(\frac{\gamma^{H}}{1 - \gamma} + \frac{1}{1 - \gamma}\sqrt{\frac{\log(2|C|/\delta)}{2m}}\right).$$
(7)

KW: $|\mathcal{C}| \leq d(d+1)/2$, hence to make second term $\leq 2\varepsilon$, enough if

$$H \geq H_{\gamma,arepsilon/\sqrt{d}} \qquad ext{and} \qquad m \geq rac{d}{(1-\gamma)^2arepsilon^2}\lograc{d(d+1)}{\delta}$$

Total # samples: $|\mathcal{C}|Hm \approx \frac{d^3 \log(d/\varepsilon) \log(d/\delta)}{(1-\gamma)^3 \varepsilon^2}$

Note: Can use $|\mathcal{C}| = \tilde{O}(d)$, paying factor of two blow-up factor: $d^3 \Rightarrow d^2$

Progress Lemma with Approximation Errors

Lemma (Geometric progress lemma with approximate policy improvement): Consider a memoryless policy π and its corresponding value function v^{π} . Let π' be any policy and define $\varepsilon : S \to \mathbb{R}$ via

$$Tv^{\pi}=T_{\pi'}v^{\pi}+arepsilon$$
 .

Then,

$$\|v^*-v^{\pi'}\|_\infty \leq \gamma \|v^*-v^\pi\|_\infty + rac{1}{1-\gamma}\,\|arepsilon\|_\infty.$$

Approximate Policy Iteration

$$Tv^{\pi_k} = T_{\pi_{k+1}}v^{\pi_k} + arepsilon_k$$

Theorem (Approximate Policy Iteration): Let $(\pi_k)_{k\geq 0}$, $(\varepsilon_k)_k$ be such that (11) holds for all $k \geq 0$. Then, for any $k \geq 1$,

$$\|v^*-v^{\pi_k}\|_\infty \leq rac{\gamma^k}{1-\gamma} + rac{1}{(1-\gamma)^2} \max_{0\leq s\leq k-1} \|arepsilon_s\|_\infty \,.$$

Approximate Policy Iteration 2

 $q_k=q^{\pi_k}+arepsilon_k',\qquad M_{\pi_k}q_k=Mq_k\,,\quad k=0,1,\ldots\,.$

Corollary (Approximate Policy Iteration with Approximate Action-value Functions): The sequence defined in (13) is such that

$$\|v^*-v^{\pi_k}\|_\infty \leq rac{\gamma^k}{1-\gamma} + rac{2}{(1-\gamma)^2} \max_{0\leq s\leq k-1} \|arepsilon_s'\|_\infty\,.$$

Least-squares Policy Iteration with G-optimal design (MC-LSPI)

- 1 Given the feature map φ , find ${\mathcal C}$ and ρ as in the Kiefer–Wolfowitz theorem
- ² Let $\theta_{-1} = 0$
- 3 For $k=0,1,2,\ldots,K-1$ do
- 4 Roll out with policy $\pi := \pi_k$ for H steps to get the targets $\hat{R}_m(z)$ where $z \in \mathcal{C}$ and $\pi_k(s) = rg\max_a \langle heta_{k-1}, \varphi(s, a)
 angle$
- 5 Solve the weighted least-squares problem given by Eq. (4) to get θ_k .
- 6 Return θ_{K-1}

$$\hat{ heta} = rg\min_{ heta \in \mathbb{R}^d} \sum_{z \in \mathcal{C}} arrho(z) \Big(\langle heta, arphi(z)
angle - \hat{R}_m(z) \Big)^2$$
 (4)

$$\hat{R}_m(z) = rac{1}{m}\sum_{j=1}^m\sum_{t=0}^{H-1}\gamma^t r_{A_t^{(j)}}(S_t^{(j)})$$

Theorem (LSPI performance): Fix an arbitrary full rank feature-map $\varphi : S \times A \to \mathbb{R}^d$ and let $K, m, H \ge 1$. Assume that $B_{2_{\varepsilon}}$ holds. Then, for any $0 \le \zeta \le 1$, with probability at least $1 - \zeta$, the policy π_K which is greedy with respect to $\Phi \theta_{K-1}$ is δ -suboptimal with



In particular, for any $\varepsilon' > 0$, choosing K, H, m so that

$$egin{aligned} &K \geq H_{\gamma,\gammaarepsilon'/2} \ &H \geq H_{\gamma,(1-\gamma)^2arepsilon'/(8\sqrt{d})} & ext{and} \ &m \geq rac{32d}{(1-\gamma)^4(arepsilon')^2} \mathrm{log}((d+1)^2K/\zeta) \end{aligned}$$

policy π_K is δ -optimal with

$$\delta \leq rac{2(1+\sqrt{d})}{(1-\gamma)^2}\,arepsilon+arepsilon'\,,$$

while the total computation cost is $\operatorname{poly}(\frac{1}{1-\gamma}, d, A, \frac{1}{(\varepsilon')^2}, \log(1/\zeta)).$

Offline planning

- This was called "global planning"
- Idea: Use the simulator to get a policy. Then keep the policy, and run with it.
- No simulator needed while using the policy.
- LSPI is offline planner with global access

From global to local access

https://arxiv.org/abs/2108.05533 Efficient Local Planning with Linear Function Approximation Dong Yin, Botao Hao, Yasin Abbasi-Yadkori, Nevena Lazić, Csaba Szepesvári

Confident MC-LSPI

Algorithm 2 CONFIDENT MC-LSPI

- 1: Input: initial state ρ , initial policy parameter w_0 , number of iterations K, regularization coefficient λ , threshold τ , discount γ , number of rollouts m, length of rollout n.
- 2: $C \leftarrow \emptyset$ // Initialize core set.
- 3: for $a \in \mathcal{A}$ do
- 4: **if** $C = \emptyset$ or $\phi(\rho, a)^{\top} (\Phi_{C}^{\top} \Phi_{C} + \lambda I)^{-1} \phi(\rho, a) > \tau$ then
- 5: $\mathcal{C} \leftarrow \mathcal{C} \cup \{(\rho, a, \phi(\rho, a), \mathsf{none})\}$
- 6: **end if**
- 7: end for
- 8: $z_q \leftarrow$ none, $\forall z \in C$ // Policy iteration starts.
- 9: for k = 1, ..., K do
- 10: for $z \in C$ do
- 11: status, result \leftarrow CONFIDENTROLLOUT $(m, n, w_{k-1}, \gamma, z_s, z_a, \Phi_{\mathcal{C}}, \lambda, \tau)$
- 12: **if** status = done, **then** $z_q \leftarrow$ result; **else** $\mathcal{C} \leftarrow \mathcal{C} \cup \{\text{result}\}$ and **goto** line (*)
- 13: **end for**
- 14: $w_k \leftarrow (\Phi_{\mathcal{C}}^\top \Phi_{\mathcal{C}} + \lambda I)^{-1} \Phi_{\mathcal{C}}^\top q_{\mathcal{C}}$
- 15: end for
- 16: return w_{K-1} , or equivalently, the policy $\pi_{w_{K-1}}$ in the form of Eq. (3.1).

(*)

Algorithm 1 CONFIDENTROLLOUT

1: Input: number of rollouts m, length of rollout n, policy parameter w, discount γ , initial state s_0 , initial action a_0 , feature matrix Φ_c , regularization coefficient λ , threshold τ .

 $\boldsymbol{\omega}$

- 2: for i = 1, ..., m do
- 3: $s_{i,0} \leftarrow s_0, a_{i,0} \leftarrow a_0$, query the simulator, obtain reward $r_{i,0} \leftarrow r(s_{i,0}, a_{i,0})$, next state $s_{i,1}$.

✓ 1

- 4: **for** t = 1, ..., n **do**
- 5: Query the simulator with $(s_{i,t}, a), \forall a \in \mathcal{A}$.
- 6: Obtain rewards $r(s_{i,t}, a)$, feature vectors $\phi(s_{i,t}, a)$, next states $s'_{i,t}(a)$, $\forall a \in \mathcal{A}$.
- 7: if there exists $a' \in \mathcal{A}$ such that $\phi(s_{i,t}, a')^{\top} (\Phi_{\mathcal{C}}^{\top} \Phi_{\mathcal{C}} + \lambda I)^{-1} \phi(s_{i,t}, a') > \tau$ then
- 8: status \leftarrow uncertain, result $\leftarrow (s_{i,t}, a', \phi(s_{i,t}, a'), \text{none})$
- 9: **return** status, result
- 10: **end if**

11:
$$a_{i,t} \leftarrow \arg \max_{a \in \mathcal{A}} w^{\top} \phi(s_{i,t}, a), r_{i,t} \leftarrow r(s_{i,t}, a_{i,t}), s_{i,t+1} \leftarrow s'_{i,t}(a_{i,t}).$$

- 12: **end for**
- 13: end for
- 14: status \leftarrow done, result $\leftarrow \frac{1}{m} \sum_{i=1}^{m} \sum_{t=0}^{n} \gamma^{t} r_{i,t}$
- 15: return status, result

<u>Theorem</u>

Under uniform ζ realizability, w.p. $\geq 1 - \delta$, Confident MC-LSPI obtains an

 $\sqrt{d} \zeta H^2 + \varepsilon$

optimal policy with at most

$$O\left(\operatorname{poly}(d, A, B, H, 1/\varepsilon, \log(1/\delta))\right)$$

simulation calls, while the total computation cost is also polynomial in the same factors.

Here *B* is a bound on the 2-norm of the parameter vectors of policies.

Can we do better?

Theorem (Query lower bound: large action sets): For any $\varepsilon > 0$, $0 < \delta \le 1/2$, positive integer d and for any (δ, ε) -sound online planner \mathcal{P} there exists a "featurized-MDP" (M, φ) with rewards in [0, 1] with $\varepsilon^*(M, \Phi) \le \varepsilon$ such that when interacting with a simulator of (M, φ) , the expected number of queries used by \mathcal{P} is at least

$$\Omega\left(\exp\left(rac{1}{32}\left(rac{\sqrt{d}arepsilon}{\delta}
ight)^2
ight)
ight)$$

.. just another needle in the haystack argument..

and when A = O(1)?

Theorem (Query lower bound: small action sets, fixed-horizon objective): For $\varepsilon > 0$, $0 < \delta \le 1/2$ and positive integer d, let

$$u(d,arepsilon,\delta) = \left\lfloor \exp\left(rac{d(rac{arepsilon}{2\delta})^2}{8}
ight)
ight
floor$$
 .

Then, for any $\varepsilon > 0, 0 < \delta \le 1/2$, positive integers A, H, d such that $d \le A^H$ and for any online planner \mathcal{P} that is (δ, ε) -sound for MDPs with at most A actions and the H-step criterion, there exists a "featurized-MDP" (M, φ) with A actions and rewards in [0, 1] such that when interacting with a simulator of (M, φ) , the expected number of queries used by \mathcal{P} is at least

$$ilde \Omega\left(rac{u(d,arepsilon,\delta)}{d(arepsilon/\delta)^2}
ight)$$

provided that $\mathrm{A}^H > u(d,arepsilon,\delta)$ ("large horizons"), while it is

$$\tilde{\Omega}\left(\frac{\mathbf{A}^{H}}{H}\right)$$

otherwise ("small horizon").



Questions from slack

Ehsan Imani 3 days ago

ML textbooks usually motivate PCA by noting that oftentimes real-world data is mostly within a small linear subspace of our d-dimensional space. Would an assumption like this rule out the pathological case of nearly-orthogonal vectors that lead to the lower bound?

+9

Jiamin He 11 hours ago

The approximation error appears in the bound of LSPI,

 $2(1+\sqrt{d})\varepsilon/(1-\gamma)^2,$

can not be controlled by the algorithm. Looks like it is also a major concern of some people according to the endnote. Why shouldn't we be worried about it? Is it possible that controlling the size of the ball (epsilon) may also increase the dimension d? Would a large value of γ be a problem here?

 $+\infty$?