

Hello!

CS292F StatRL Lecture 9

Exploration in Tabular MDPs

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara


Logistic notes

- HW1 due today.
- HW2 is posted on the course website.
 - Q1: A simple coding question
 - Q2: An alternative rate-optimal algorithm for MAB.
 - Q3: Exploration in tabular RL

Recap: Lecture 8

- Linear Bandits
 - Problem setup
 - Regret definition
- Optimism in the face of uncertainty
 - LinUCB algorithm
 - Bounding sum of square regrets with information gain.
 - a self-normalized Martingale Concentration

This lecture: Exploration in Reinforcement Learning

- Why is it challenging?
 - The reward depends on both s , a
 - Unlike the generative model setting, we cannot just choose any s to explore.
 - The data needs to be actively collected
- We will study
 - Tabular MDP 
 - Linear MDPs
 - Both in the finite horizon episodic setting.

Recap: Finite horizon MDPs

- Parameterization / Setup

$$M = (\mathcal{S}, \mathcal{A}, \{P\}_h, \{r\}_h, H, \mu)$$

O-based index.
 P_0, P_1, \dots, P_{H-1}

horizon $H < +\infty$

$S_0 \cup \mathcal{M}$

$\mu = S_0$

Nonstationary Transition: $P_h(\cdot | S, A) \rightarrow \Delta(\mathcal{S})$ differs for each h

r_0, \dots, r_{H-1}

$R_h: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

$E[R_h(S_h=S, A_h=a)] =: V_h(S, a)$

- Additional notations

- Q functions Q_h^π, Q_h^* $h=0, 1, \dots, H-1$

- V functions V_h^π, V_h^*

- Policies $\pi = \{ \pi_0, \dots, \pi_{H-1} \}$

$$Q_H^\pi = Q_H^* = 0 \quad \forall s, a \text{ (base)}$$

for convergence

- Observed trajectory data

Execute policy π :

trajectory $\rightarrow \left\{ (S_0^k, A_0^k), (S_1^k, A_1^k), \dots, (S_{H-1}^k, A_{H-1}^k) \right\}$

$$Q_h^\pi(s, a) = V_h(s, a) + E^\pi_{(S', A') \sim P_h^\pi(S', A' | s, a)} Q_{h+1}^\pi(S', A')$$

all $h=0, 1, \dots, H-1$

Problem setup: online learning of Finite horizon MDPs

- Agent decides on a policy

π_0 at 0th Episode, \dots , π_k at k th Episode

$$\pi_k = \{ \pi_{k,h} \text{ for } h=0, \dots, H \}$$

$$\pi_{k,h}: S \rightarrow \Delta(A)$$

- Collect a trajectory

$(S_0, A_0) \dots (S_{H-1}, A_{H-1})$, ~~(S_H, A_H)~~

- Agent updates the policy.

$$\pi_{k+1} \leftarrow \text{function}(\text{Hist}_k)$$

- Regret definition

$$\text{Regret} = \sum_{k=0}^{K-1} \text{Regret}_k = \underbrace{K \cdot V^*(s_0)}_{\uparrow} - \sum_{k=0}^{K-1} V^{\pi_k}(s_0)$$

$$E[\text{Regret}] = K V^*(s_0) - E\left[\sum_{k=0}^{K-1} V^{\pi_k}(s_0)\right]$$

$$\pi_h^*(s) = \underset{a \in A}{\text{arg max}} Q_h^*(s, a)$$

for each h

Recap: The need for strategic exploration

$$r(S_i, a) = 0 \quad \forall i < H-1 \quad \forall a$$

$$r(S_{H-1}, a) = 1 \quad \forall a$$



Randomized exploration

$$P(S_{H-1}) = d \left(\frac{1}{2} \right)^H$$

UCB-VI: model-based learning by **optimistic** value Iterations

- Construct estimates of the transition kernels

$$\hat{P}_h^k$$

- Design exploration bonuses

At Episode k : $\hat{b}_h^k(s,a)$

$Q_h(s,a) + \hat{b}_h^k(s,a) \geq Q_h^*(s,a)$

model-based approach

$$\hat{M}^k$$

- Idea: based on the uncertainty in the transition kernel estimates

- Update the policy by **optimistic** value iteration

How do we estimate the model parameters (P and r)?

- Simple plug-in estimator

At $k, h,$

$$\hat{P}_h^k(S'|S_a) = \frac{N_h^k(S, a, S')}{N_h^k(S, a)}$$

$N_h^k(S, a, S')$ = # of times these triplets appear from step h to step $h+k-1$

$$= \sum_{i=0}^{k-1} \mathbb{I}(S_h^i = S, A_h^i = a, S_{h+i}^i = S')$$

$$N_h^k(S, a) = \sum_{i=0}^{k-1} \mathbb{I}(S_h^i = S, A_h^i = a)$$

- What happens if we observe no state-action pairs?

$$\frac{0}{0} =: 0$$

What does value iteration do in finite horizon MDPs?

$$0 \leq \Pr_h(s,a) \leq 1$$

$$\Downarrow$$

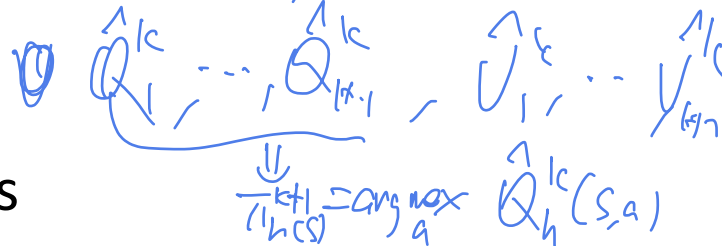
$$\hat{Q}_h^k(s,a) \leq H$$

$$\hat{V}_H^k(s) = 0, \forall s,$$

$$\hat{Q}_h^k(s, a) = \min \left\{ r_h(s, a) + b_h^k(s, a) + \hat{P}_h^k(\cdot | s, a) \cdot \hat{V}_{h+1}^k, H \right\},$$

$$\hat{V}_h^k(s) = \max_a \hat{Q}_h^k(s, a), \pi_h^k(s) = \operatorname{argmax}_a \hat{Q}_h^k(s, a), \forall h, s, a.$$

for $h = H, \dots, 0$



- Remark:

- It converges in H steps
- It produces a non-stationary policy indexed by h

How do we design exploration bonuses?

scaffolding style bonus

$$b_h^k(s, a) = \underbrace{H}_{\substack{\uparrow \\ \text{max possible } V \text{ function}}} \sqrt{\frac{L}{N_h^k(s, a)}} \quad \text{where } \underline{L := \ln(SAHK/\delta)}$$

- Intuitively, this encourages exploring new state-action pairs.
- Idea: propagate errors from the estimated transitions over to the rewards.

The regret of UCB-VI (Azar et al., 2017) (Jackisch et al., 2019)

- Theorem (AJKS Thm 6.1):

$$\text{Regret} := \mathbb{E} \left[\sum_{k=0}^{K-1} \left(V^* - V^{\pi^k} \right) \right] \leq 2H^2 S \sqrt{AK \cdot \ln(SAH^2 K^2)} = \tilde{O} \left(H^2 S \sqrt{AK} \right)$$

$$\tilde{O}(\sqrt{H^4 S^2 AK})$$

N₂ regret learning Alg

- This is not optimal in H, S, but a simple analysis to start. We will talk about how to improve it towards the end.

$$\text{lower bound: } \Omega(\sqrt{H^3 S AK})$$

nonstationary transition
MDP
H + 200
 12

Step 1: Concentration

Condition on $N_h^k(s_a)$
 $\forall h, k$ for $i=0, 1, \dots, k-1$ are iid

Lemma 1: for all h, k, s, a , w.p. $1-\delta$

$$\left\| \hat{P}_h^k(\cdot | s_a) - P_h^k(\cdot | s_a) \right\|_1 \leq \left(\frac{S \log \frac{S^2 H K}{\delta}}{N_h^k(s_a)} \right)$$

By McDiarmid Inequality

Lemma 2: w.p. $\geq 1-\delta$, for $\forall s, a, h, k$

$$\left\| \hat{P}_h^k(\cdot | s_a) \cdot V_{h+1}^* - P_h^k(\cdot | s_a) \cdot V_{h+1}^* \right\| \leq H \sqrt{\frac{L}{N_h^k(s_a)}}$$

$$\frac{1}{N_h^k(s_a)} \sum_{i=1}^{k-1} \underbrace{V_{h+1}^*(S_{h+1}^i)}_{\text{iid R.V.}} \mathbb{1}(S_h^i = s, A_h^i = a)$$

$$\underline{P(\text{Fail})} \leq 2\delta$$

Step 2: Optimism

under "not fail"

Lemma 3, $\hat{V}_h^k \geq V_h^*$ for all $h=0, 1, \dots, H-1, H$

Proof: Base: $\hat{V}_H^k = V_H^* = 0$

Assume for h , $\hat{V}_h^k \geq V_h^*$, we will prove that $\hat{V}_{h-1}^k \geq V_{h-1}^*$

$$\hat{V}_{h-1}^k = \max_a Q_{h-1}^k(\cdot, a), \text{ and } Q_{h-1}^k(s, a) = \min_{h'} \left(b_{h-1}^k(s, a) + \sum_{h'} P_{h'}^k(s, a) \cdot \hat{V}_{h'}^k \right)$$

Case ① when H is the smaller: $Q_{h-1}^k(s, a) = H \geq Q_{h-1}^*(s, a)$

Case ② when H is not selected

$$Q_{h-1}^*(s, a) = V_{h-1}^*(s, a) + \sum_{h'} P_{h'}^*(s, a) \cdot V_{h'}^*(s, a)$$

show $\hat{V}_{h-1}^k(s, a) \geq V_{h-1}^*(s, a)$

$$\begin{aligned} \hat{V}_h^k(s) &= \max_a \hat{Q}_h^k(s, a) \\ &\geq \hat{Q}_h^k(s, a^*) \geq Q_h^*(s, a^*) \\ &= V_h^*(s) \quad \square \end{aligned}$$

$$\begin{aligned} \hat{Q}_{h-1}^k(s, a) - Q_{h-1}^*(s, a) &= b_{h-1}^k(s, a) + \sum_{h'} P_{h'}^k(s, a) \hat{V}_{h'}^k - \sum_{h'} P_{h'}^*(s, a) V_{h'}^* \\ &\geq b_{h-1}^k(s, a) + \left(\sum_{h'} P_{h'}^k(s, a) - \sum_{h'} P_{h'}^*(s, a) \right) \cdot V_h^* \end{aligned}$$

$$\geq \underbrace{b_{h-1}^k(s, a) - H \cdot \frac{2}{\sqrt{H}}}_{\geq 0}$$

Finite horizon simulation lemma

(from HW1) Q5.

$$\boxed{\| \hat{V}_h^{\pi} - V_h^{\pi} \| \leq \text{small}}$$

$$\hat{V}_h^{\pi} = r_h + P_h^{\pi} \hat{V}_{h+1}^{\pi}$$

$$V_h^{\pi} = r_h + P_h^{\pi} V_{h+1}^{\pi}$$

Take difference

$$\hat{V}_h^{\pi} - V_h^{\pi} = \hat{r}_h - r_h + \underbrace{P_h^{\pi} (\hat{V}_{h+1}^{\pi} - V_{h+1}^{\pi})}_{\text{add}} + \underbrace{(P_h^{\hat{\pi}} - P_h^{\pi}) \hat{V}_{h+1}^{\hat{\pi}}}_{\text{subtract}} \quad \text{(*)}$$

recursively apply the difference

$$\hat{V}_0^{\hat{\pi}} - V_0^{\pi} = \hat{r}_0 - r_0 + P_0^{\hat{\pi}} (\hat{V}_1^{\hat{\pi}} - V_1^{\pi}) + (P_0^{\hat{\pi}} - P_0^{\pi}) \hat{V}_1^{\hat{\pi}}$$

$$P_0^{\hat{\pi}} \left[\hat{r}_1 - r_1 + (P_1^{\hat{\pi}} - P_1^{\pi}) \hat{V}_1^{\hat{\pi}} + P_1^{\hat{\pi}} (\hat{V}_2^{\hat{\pi}} - V_2^{\pi}) \right]$$

$$\begin{aligned} \hat{V}_0^{\hat{\pi}} - V_0^{\pi} &= \sum_{h=0}^{H-1} \mathbb{E}^{\pi} \left[\hat{r}_h(s_h) - r_h(s_h) + (P_h^{\hat{\pi}}(\cdot | s_h) - P_h^{\pi}(\cdot | s_h)) \cdot \hat{V}_{h+1}^{\hat{\pi}} \right] \\ &= \sum_{h=0}^{H-1} \mathbb{E}^{\pi} \left[\sum_a \pi_h(a | s_h) (\hat{r}_h(s_{h+1}, a) - r_h(s_{h+1}, a)) + \hat{P}_h^{\hat{\pi}}(\cdot | s_{h+1}, a) - P_h^{\pi}(\cdot | s_{h+1}, a) \cdot \hat{V}_{h+1}^{\hat{\pi}} \right] \end{aligned}$$

Regret in kth Episode

Simulation lemma and also take $\min(H, \frac{1}{\epsilon})$
 make it smaller

$$\begin{aligned}
 \text{Regret}_k &= V_0^*(s_0) - V_0^{\text{Tric}}(s_0) \stackrel{\text{By optimality}}{\leq} V_0^{\text{Tric}}(s_0) - V_0^{\text{Tric}_k}(s_0) \stackrel{\text{Simulation lemma}}{\leq} \sum_{h=0}^{H-1} E^{\text{Tric}} \left[r_h(s_h, A_h) - v(s_h, A_h) \right. \\
 &\quad \left. + (P_h^{\text{Tric}}(\cdot | s_h, A_h) - P_h^{\text{Tric}_k}(\cdot | s_h, A_h)) \cdot V_{h+1}^{\text{Tric}} \right] \\
 &= \sum_{h=0}^{H-1} E^{\text{Tric}} \left[b_h^k(s_h, A_h) + (P_h^{\text{Tric}}(\cdot | s_h, A_h) - P_h^{\text{Tric}_k}(\cdot | s_h, A_h)) \cdot V_{h+1}^{\text{Tric}} \right] \\
 &\leq \sum_{h=0}^{H-1} E^{\text{Tric}} \left[\geq H \sqrt{\frac{SL}{N_h^{\text{Tric}}(s_h, A_h)}} \right] \leq \|P_h^{\text{Tric}}(\cdot | s_h, A_h) - P_h^{\text{Tric}_k}(\cdot | s_h, A_h)\|_1 \|V_{h+1}^{\text{Tric}}\|_{\infty} \\
 &= \underbrace{2H\sqrt{SL}}_{\text{Lemma 1}} \cdot E \left[\sum_{h=0}^{H-1} \sqrt{\frac{1}{N_h^{\text{Tric}_k}(s_h, A_h)} \mid \text{hist}_k} \right] \leq \sqrt{SL} \cdot H
 \end{aligned}$$

Total regret

$$E \left[\sum_{k=0}^{K-1} \text{Regret}_k \right] = E \left[\sum_{k=0}^{K-1} V_{S_k}^* - V_{S_k}^{T_k} \right] = E \left[\mathbb{I}(\text{Not Fail}) \sum_{k=0}^{K-1} V^* - V^{T_k} \right]$$

$$+ E \left[\mathbb{I}(\text{Fail}) \sum_{k=1}^{K-1} (V^* - V^{T_k}) \right]$$

$$\leq E \left[\left(\sum_{k=0}^{K-1} V^* - V^{T_k} \right) \mathbb{I}(\text{Not Fail}) \right] + 2\delta \cdot K \cdot H$$

$$\leq E \left[\sum_{k=0}^{K-1} \sum_{h=0}^{H-1} \frac{H \sqrt{SL}}{\sqrt{N_h^{(k)}(S_h, A_h)}} \right] + 2\delta \cdot K \cdot H \quad (\Delta)$$

\swarrow Regret_k bound holds under "Not Fail"
 \nwarrow From last lecture.

Choose $\delta = \frac{1}{K \cdot H}$

$$(\Delta) = \sum_{h=0}^{H-1} \sum_{S, A \in \mathcal{S} \times \mathcal{A}} \sum_{i=1}^{N_h^{(k)}(S, A)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{S, A} 2 \sqrt{N_h^{(k)}(S, A)}$$

Regret $\leq 4H^2 S \sqrt{AK}$

$$(\Delta) \leq E \left[\sum_{h=0}^{H-1} \sum_{S, A} 2 \sqrt{N_h^{(k)}(S, A)} \right] + 2\delta \cdot K \cdot H$$

$$\leq \sum_{h=0}^{H-1} \sum_{S, A} \sqrt{SA \cdot \sum_{i=1}^{N_h^{(k)}(S, A)} \frac{1}{i}} + 2\delta \cdot K \cdot H$$

$$\leq 2H \sqrt{SA \cdot K} + 2\delta \cdot K \cdot H \leq 4H \sqrt{SAK}$$

Ideas for improving the dependence on S and H

$$\Omega(\sqrt{H^3 S A K})$$

$$O(\sqrt{H^4 S^2 A K})$$

Improve S: $\left(\underbrace{\hat{P}_n^{\text{loc}}(S_n, A_n)}_{\text{Lemma 2}} - \underbrace{\hat{P}_n^{\text{true}}(S_n, A_n)}_{\downarrow} \right) \left(\underbrace{\hat{V}_{n+1}^{\text{loc}} - V_{n+1}^*}_{\downarrow \frac{1}{\sqrt{N}}} + V_{n+1}^* \right)$

$$= \underbrace{(\hat{P} - P) \cdot V^*}_{\text{Lemma 2}} + \underbrace{\sqrt{\frac{SL}{N}} \cdot \frac{H \sqrt{SL}}{N}}_{\downarrow \frac{1}{\sqrt{N}}} = \frac{H \sqrt{L}}{N} + \frac{S H L}{N}$$

Improve H: Bernstein's inequality

Final notes about exploration in Tabular MDPs

- Optimal rates:

- Non-stationary transitions
- Stationary transitions

$$\begin{aligned} & \Omega(\sqrt{H^3 S A K}) \\ & \Omega(\sqrt{H^2 S A K}) \end{aligned}$$

- State of the art:

- Stationary case: MVP $O(\sqrt{H^2 S A K} + H^2 S^2 A)$
 - Zhang, Ji and Du (2020) <https://arxiv.org/pdf/2009.13503.pdf>
 - Modified the episode reward bound from $[0,1]$ to $[0,H]$ to be consistent with this lecture

- Nonstationary case: $O(\sqrt{H^3 S A K} + H^4 S^2 A)$

- Q-learning: [Jin et al.](#), [Bai et al.](#), optimal rates in [Zhang et al. \(2020\)](#)

- Open problem:

- Is it possible to get rid of the S dependence in the low-order terms.