

Lecture 21

CS 59006 - RL

MDPs

Agenda

- VCR L2 (Pebe, Josch, Ronald 2010)

$$\text{Aim: } \hat{R}_T = \underline{T} - \sum_{t=1}^T r(X_t, A_t)$$

what is oracle
long term return

what our algorithm
achieves

- $\bar{r}(x, a) = r(x, a)$ deterministic, known

- Diameter $D < \infty$.

$$-0 \leq r(x, a) \leq 1$$

$$T_b(x, a) = \sum_{i=1}^t \mathbb{I} \{X_i = x, A_i = a\}$$

$$\rightarrow P_b(x' | x, a) = \frac{\sum_{i=1}^t \mathbb{I} \{X_i = x, A_i = a, X_{i+1} = x'\}}{\max_{x'} T_t(x, a)}$$

Side note

- binomial: $p \rightarrow$ w.p. $p \rightarrow \text{outcome } \begin{bmatrix} 1 \\ 0 \end{bmatrix}$
w.p. $1-p \rightarrow \text{outcome } \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

- multinomial distribution $p \in \Delta^d \rightarrow \text{outcome } e_i \text{ w.p. } p_i$

→ $P \in \Delta^d$ means $P \in \mathbb{R}_+^d$ and $\sum_{i=1}^d p_i = 1$
 ↑
 simplex.

This is a die with d sides, side i is the outcome with probability P_i

If we roll this die n times
 what is \hat{P} ?

$$\frac{1}{n} \sum_{t=1}^n e_{x_t} = \hat{P} \rightarrow \hat{P}_i = \frac{\sum_{t=1}^n \mathbb{I}(x_t = i)}{n}$$

we hope $\hat{P} \approx P$

$$e_i = \begin{bmatrix} \vdots \\ i \\ \vdots \end{bmatrix} \leftarrow 1 \text{ at } i\text{th index}$$

Theorem (Concentration on simplex)

[Weissman et al. 2003]

Let X_1, \dots, X_n be i.i.d sequence of random

variables on $(\Omega, \mathcal{F}, \mathbb{P})$ such that $X_t: \Omega \rightarrow \{1, \dots, m\}$

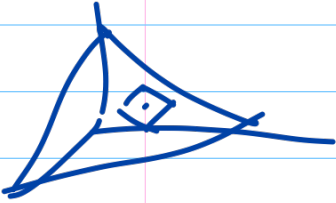
Let $p_i = \mathbb{P}(X_1 = i)$ and $\hat{P}_i = \frac{1}{n} \sum_{t=1}^n \mathbb{I}\{X_t = i\}$

Then:

$$\mathbb{P} \left(\|P - \hat{P}\|_1 > \sqrt{\frac{2m \log \frac{2}{\delta}}{n}} \right) \leq \delta$$

Now back to MDPs and UCR2:

$$\rightarrow C_t^\delta(x, a) = \{P' \in \Delta^{|X|} : \|P' - P_{t-1}(\cdot | x, a)\|_1,$$



$$\leq \sqrt{\frac{|X| L_{t-1}(x, a)}{\max(1, T_{t-1}(x, a))}}$$

$$\text{for } L_{t-1}(x, a) = 2 \log \left(\frac{4 |X| |A| T_t(x, a) (1 + T_b(x, a))}{\delta} \right)$$

(If $T_t(x, a) = 0 \Rightarrow \text{set } L_{t-1}(x, a) = 1$)

$$\rightarrow C_b^\delta = \{P' = \{P'(\cdot | x, a)\}; P'(\cdot | x, a) \in C_t^\delta(x, a), \forall (x, a) \in X \times A\}$$

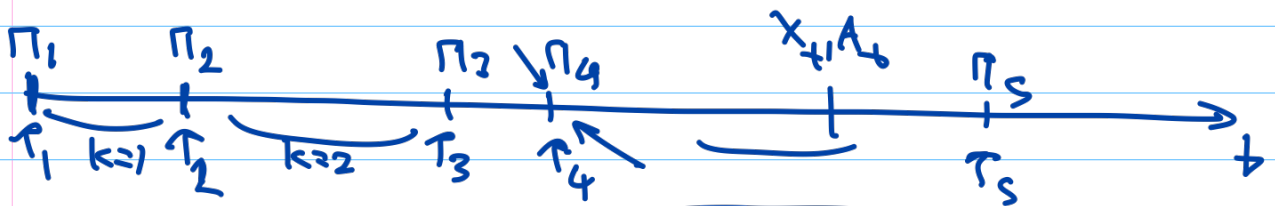
↑ The set of plausible model.

$$\text{For } L_t(x, a) \leq L := 2 \log \left(\frac{4 |X| |A| T(T+1)}{\delta} \right)$$

We want the true probability kernel P to be in C_t^δ for all T , with probability at least $1 - \delta$.

UCRL 2:

- start with a policy
- collect samples
- stop when the number samples at least for a pair (n, a) is doubled,
- compute $C_b^{\delta}()$ and P_t
- compute optimistic model \tilde{P}_t
- Compute π_k , the optimistic policy
- Then start new epoch $k+1$



If $T_t(X_t, A_t) < 2T_{T_{k-1}}(X_t, A_t)$, you continue
otherwise stop the epoch.

How to come up with optimistic model.

we want

$$\rho_k = \max_{a \in \mathcal{A}} \max_{\pi} \max_{P' \in \mathcal{C}_{\rho_k}} \rho_{\pi}^{\pi}(P')$$

The best \tilde{P}_k

and $\rho_k + V_k(n) = \max_{a \in \mathcal{A}} r(n, a)$

$$+ \langle \tilde{P}_k(\cdot | n, a), V_k \rangle$$

$\forall (n, a) \in \mathcal{X} \times \mathcal{A}$

Define \tilde{M}_k , the extended MDP with the same state space as M , but extended action space

$$\tilde{A}_k = \{(a, P') : a \in A, P' \in C_k^n(n, a)\}$$

with the reward $r(n, (a, P')) = r(n, a)$

and transition $P(\cdot | n, (a, P')) = P'(\cdot | n, a)$

For any policy Π in M , there exists $\tilde{\Pi}$ on \tilde{M}

such that for any $n \rightarrow a = \Pi(n)$ on M
 $\Pi(n), P = \tilde{\Pi}(n)$ on \tilde{M}

$$\Rightarrow D(\tilde{M}) \leq D$$

Solve for P and V at \tilde{M} , then Π_k follows

$$\Pi_k(n) \in \arg \max_{a \in A} r(n, a) + \max_{P' \in C_k^n(n, a)} \langle P'(n, a), V \rangle$$

\downarrow \downarrow
 $P'(n, a) \in C_k^n(n, a)$

There is an efficient alg to solve it (UCRL2)

Lemma: For a (ρ, V) , satisfying Bellman optimality equation:

$$\text{span}(V) \leq D \quad \left[\text{UCRL2, Chp 38 Bandit book} \right]$$

Solving Bellman optimality for \tilde{M}_k give

$\rho_k, \tilde{\rho}_k, V_k$. Since V_k is not unique and $\text{span}(V_k) \leq D$, we choose V_k s.t. $\|V_k\|_\infty \leq \frac{D}{2}$

Proof of regret: $\hat{R} \leq CD|x| \sqrt{|A|T \log\left(\frac{T|x|W}{\delta}\right)}$

$$\hat{R}_T = \sum_{t=1}^T \left(\overset{\downarrow}{\rho} \underset{\uparrow}{r}(x_t | A_t) \right) = \sum_{k=1}^K \sum_{t=\tau_k}^{\tau_{k+1}-1} (\rho - r(x_t | A_t))$$

→ Optimism step:

$$\leq \sum_{k=1}^K \left(\sum_{t=\tau_k}^{\tau_{k+1}-1} (\rho_k - r(x_t | A_t)) \right)$$

with probability at least $1-\delta$.

→ R_k

Also we know: $\rho_k = r_{\pi_k}(x_t) - V_k(x_t) + \tilde{\rho}_k(\cdot | x_t, \pi_k(x_t), V_k)$

Monday, November 2, 2020

$$\begin{aligned}
 R_k &\leq \sum_{t=\tau_k}^{\tau_{k+1}-1} \left(-V_k(x_t) + \langle \tilde{P}_k(\cdot | x_t, A_t), V_k \rangle \right) \\
 &= \sum_{t=\tau_k}^{\tau_{k+1}-1} \left(-V_k(x_t) + \langle P(\cdot | x_t, A_t), V_k \rangle \right) \rightarrow A_1 \\
 &\quad + \sum_{t=\tau_k}^{\tau_{k+1}-1} \langle \tilde{P}_k(\cdot | x_t, A_t) - P(\cdot | x_t, A_t), V_k \rangle \rightarrow A_2
 \end{aligned}$$

$$\begin{aligned}
 (A_1) &= \sum_{t=\tau_k}^{\tau_{k+1}-1} \left(V_k(x_{t+1}) - V_k(x_t) + \langle P(\cdot | x_t, A_t), V_k \rangle - V_k(x_{t+1}) \right) \\
 &\quad \downarrow \\
 &\quad V_k(x_{\tau_{k+1}}) - V_k(x_{\tau_k}) \leq D
 \end{aligned}$$

$$(A_1) \leq D + \sum_{t=\tau_k}^{\tau_{k+1}-1} \left(\mathbb{E}[V_k(x_{t+1}) | \mathcal{F}_t] - V_k(x_t) \right)$$

→ mean zero martingale sequence, and bounded adapted to \mathcal{F}_t

using H\"older inequality

$$A_2 \leq \sum_{t=\tau_k}^{\tau_{k+1}-1} \|P_k(\cdot | x_t, A_t) - P(\cdot | x_t, A)\|_1 \|V_k\|_\infty$$

from the confidence interval $D/2$

$$(A_2) \leq \frac{D \sqrt{L|X|}}{2} \sum_{(x,a) \in \mathcal{X} \times \mathcal{A}} \frac{T_{(x)}(x,a)}{\sqrt{\max(1, T_{\tau_k-1}(x,a))}}$$

Monday, November 2, 2020

$$T_{Q,K}(n,a) = \sum_{t=p_K}^{r_{K+1}} I(X_t = n, A_t = a)$$
