

(April 8)

$$H-A \approx O(\sqrt{K})$$

$$R_K \lesssim \sum_{k=1}^K \zeta_k + \sum_{k=1}^K \sum_{h=1}^{H-1} \gamma_h^{(k)}$$

$$+ 2H C_0 \left[ \sum_{k=1}^K \sum_{h=0}^{H-1} \frac{1}{\sqrt{1 \vee N_k(S_h^{(k)}, A_h^{(k)})}} \right]$$

$\forall u \in \{0, 1, \dots\}$

IV.

$$\beta(u) \leq \frac{C_0}{\sqrt{u}}$$

II

$$\beta(0) = 1 \leq \frac{C_0}{\sqrt{1}}$$

$S, A, K, H$  big  
 $\zeta$  small

$$\bar{N} = \sum_{s, a} \sum_{k=1}^K \sum_{h=0}^{H-1} \frac{\mathbb{P}(S_h^{(k)} = s, A_h^{(k)} = a)}{\sqrt{1 \vee N_k(s, a)}}$$

$$M_K(s, a) = \sum_{h=0}^{H-1} \mathbb{P}(S_h^{(k)} = s, A_h^{(k)} = a)$$

$$N_\theta(s, a) = M_1(s, a) + \dots + M_K(s, a)$$

$$= \sum_{s,a} \sum_{k=1}^K \underbrace{M_k(s,a)}_{\sqrt{1/V(M_1(s,a) + \dots + M_K(s,a))}} \xleftarrow{f'(k)}$$

$f(k)$

$$\int_1^K \frac{f'(x)}{\sqrt{f(x)}} dx = 2 \left[ \sqrt{f(x)} \right]_1^K$$

$$2(\sqrt{f})' = \frac{1}{2} \cdot \frac{f'}{\sqrt{f}} \leq 2\sqrt{f(K)}$$

Lemma:  $\forall m_1, \dots, m_K \geq 0$

$$\sum_{k=1}^K \frac{m_k}{\sqrt{1/V(m_1 + \dots + m_k)}} \leq 2 \sqrt{m_1 + \dots + m_K}$$

Proof: telescoping

$$\begin{aligned} \overline{IV}_1 &\leq 2 \sum_{s,a} \sqrt{N_\theta(s,a)} \\ &\leq 2SA \left[ \sum_{s,a} \frac{1}{SA} \sqrt{N_\theta(s,a)} \right] \\ &\leq 2SA \sqrt{\sum_{s,a} N_\theta(s,a)/SA} \\ \text{Justins} &= 2\sqrt{SAHK} \end{aligned}$$

Theorem: wp.  $1 - O(\delta)$

$$R_K \leq 4C_{SH}\sqrt{SAHK} +$$

$\uparrow$

$$2\sqrt{S\omega(2) + 4yHKSA/\delta}$$

UCRL(2)

$$C' H \sqrt{K \log \frac{1}{\delta}} +$$

Jaksch  
Ortner  
Auer

$$C'' H \sqrt{HK \log \frac{1}{\delta}}$$

Zihau Zhang / Xiaoyang Ji / Siwei S. Du

Theorem: Total reward per episode  $\leq R$

$$R_K = O\left(R \left(\sqrt{SAK} + S^2 A\right) \text{poly}\left(\frac{SAHK}{\delta}\right)\right)$$

Theorem:  $R_K = \mathcal{O}\left(R \sqrt{SAK}\right)$

Inhomogeneous :

$$P = (P_0, \dots, P_H)$$

$$r = (r_0, \dots, r_H)$$

extra  $H$

Optimistic

$Q$ -learning is opt.

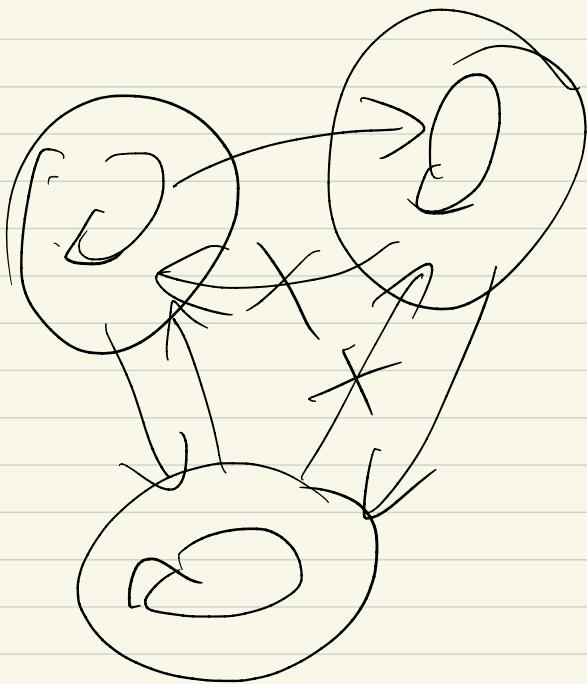
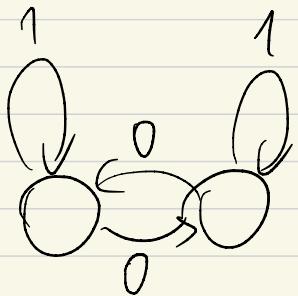
Inf. horizon

→ episodes ← end break!

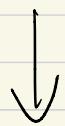
|  
count doubles

|  
fixed schedule

necessary : linear reset  
otherwise !

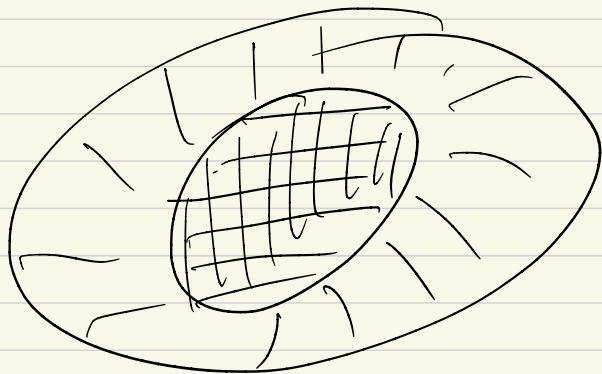


Randomize



whp you are  
optimistic

PSRL / Thompson sampling



Elimination

IDS

Finite-horizon setting:

$H > 0$

$$\ell_h: S \times A \rightarrow \mathbb{R}^d$$

$$\|\ell_h\|_2 \leq 1$$

Inhomogeneous case

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

$g^* \in \mathcal{F}_q$ ?

Idea:

(UCB)

on  $g^*$  + greedy.

$$S_0^{(k)}, A_0^{(k)}, \dots$$

$$S_H^{(k)}, A_{H-1}^{(k)}, S_H^{(k)}$$

$$q_h^{(k)} = \left[ \Phi_h W_h^{(k)} \right]^H + r_h + \beta_{bh} \quad \text{known}$$

$$v_h^{(k)} = M q_h^{(k)}$$

$$q_h^* = r_h + P_h v_{h+1}^*, \quad v_h^* = M q_h^*$$

Value iteration

Beginning of episode  $k$ ; data up to  $k-1$

$$\underline{W_h^{(k)}}(v) = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^{k-1} \underbrace{\left( q_h(Z_{ih})^T w - v \right)^2}_{+ \lambda \|w\|_2^2} + \frac{1}{2} \|w\|_2^2$$

$$Z_h = (S_h^{(k)}, A_h^{(k)})$$

$$W_h^{(k)} := \underbrace{W_h^{(k)}(v_{h+1}^{(k)})}_{+ \lambda \|w\|_2^2}$$

LSVI - UCB

$\beta_{bh}$

TBC

Chi Jin  
Zhuoran Yang Zhao Song  
M. Jordan

Ass:  $V \subseteq \mathbb{R}^{SA}$

True:  $P_h v \in \mathcal{F}_{q_h}$   
the [H]

Lemma:

$$X_1, \dots, X_{n_f} \in \mathbb{R}^d$$

$$Y_1, \dots, Y_{n_f} \in \mathbb{R}$$

$$Y_t = X_t^T \theta + \varepsilon_t$$

$$\mathbb{E}[e^{\lambda \varepsilon_t} | \mathcal{F}_{t-1}] \leq e^{\frac{\lambda^2}{2} \sigma^2}$$

$$\hat{\theta}_t = \sum_{t=1}^{\hat{T}} (X_t^T \theta - Y_t)^2 + \gamma \|\theta\|_2^2$$

$$\text{WP } 1-\delta \quad \forall x \in \mathbb{R}^d$$

$$|\langle x, \hat{\theta}_t \rangle - \langle x, \theta_* \rangle| \leq \|x\|_{M_t^{-1}} \beta_t(\delta)$$

$$M_t = \sum_{s=1}^t X_s X_s^T + \gamma I$$

$$\beta_t(\delta) = \sqrt{\log \det \left( \frac{(M_t)^{1/2}}{\gamma I} / \delta \right)} + \gamma \|\theta\|_2$$

$$v \mapsto w_n^{(k)}(v) \quad \boxed{v \mapsto \phi w_n^k(v)} = P_n^{(k)} v$$

$$P_n v$$

linear in  $v$