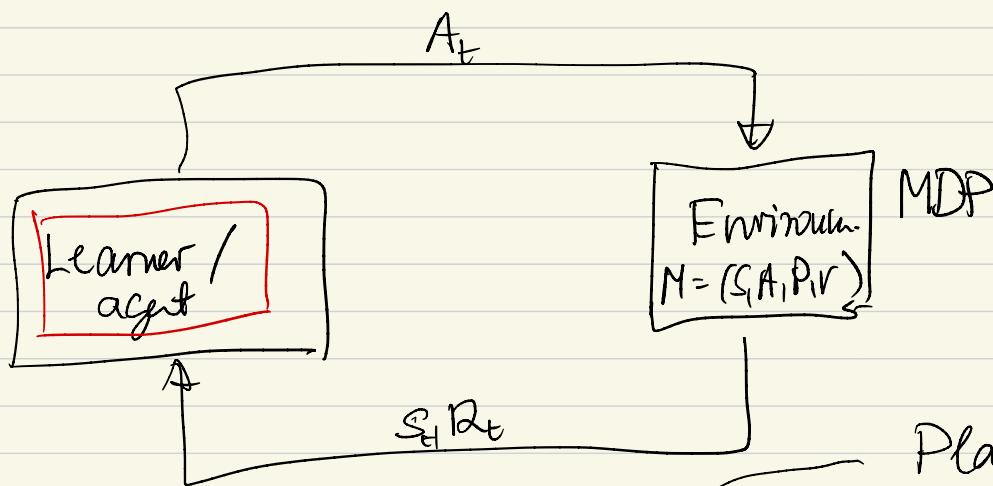
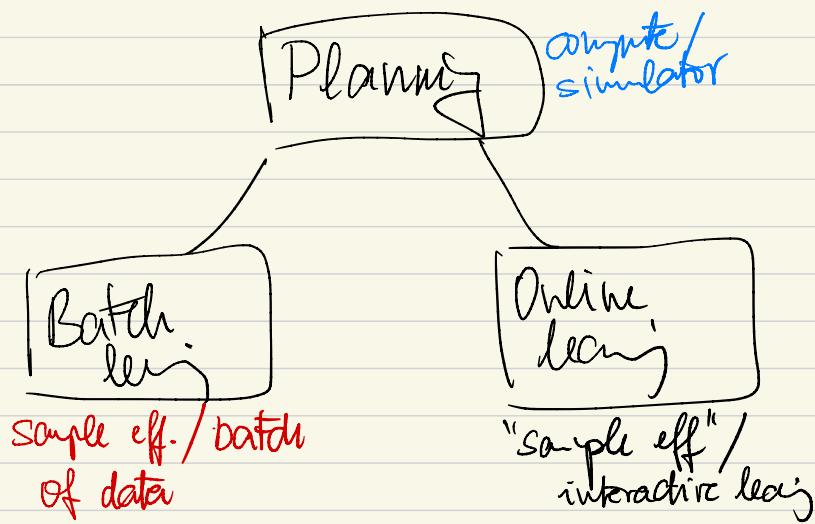


April 1

Online learning

/ Online RL



How good is a learner?

"Regret"

Cumulative regret

over time

undiscounted

fixed-horizon

infinite horizon

$$R_n = v_n^*(S_0) - \sum_{t=0}^{n-1} R_t$$

$$\rightarrow R_n = g^* n - \sum_{t=0}^{n-1} R_t$$

↑ opt. arg. val.

Episodes of length $H > 0$, $h = 0, 1, \dots, H-1$
 $1 \leq k \leq K$: # episodes.

$$S_0^{(k)} \sim \mu, A_0^{(k)}, S_1^{(k)}, A_1^{(k)}, \dots, S_{H-1}^{(k)}, A_{H-1}^{(k)}, S_H^{(k)}$$

$$S_{h+1}^{(k)} \sim P_{A_h^{(k)}}(S_h^{(k)})$$

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

For simplicity: r is known

$$R_K = \sum_{k=1}^K \left(V_0^*(S_0^{(k)}) - \sum_{h=0}^{H-1} r_{A_h^{(k)}}(S_h^{(k)}) \right)$$

\uparrow
regret

$$\mathbb{E}\left[\frac{R_K}{K}\right] \rightarrow 0, K \rightarrow \infty$$

PAC \leftrightarrow Regret

MDP

$H = 1 \Leftrightarrow$ contextual bandits

$|S| = 1$ finite-armed bandits

$$n = KH$$

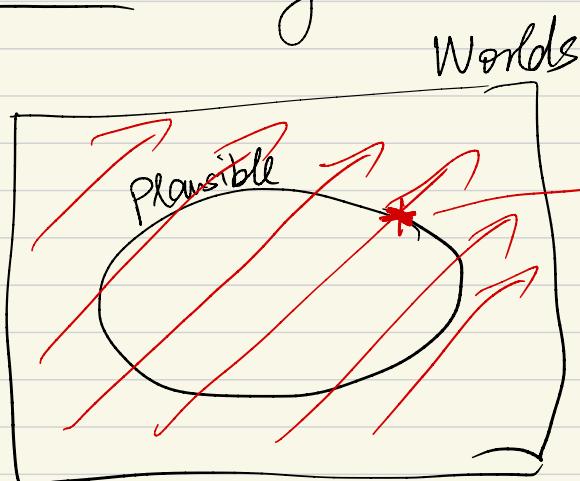
reward is unknown

$\boxed{\epsilon\text{-greedy}}$ $R_n = \Theta(n^{2/3})$, $\underline{R_n = O(\sqrt{n})}$

Principle: Optimism in the Face of Uncertainty

Necessary: Tryig s_g gives info about how good s_g is

Suff. for opt. to be "opt": $t_g x$ does not give "much" info about y



follow policy that gives best value in the best world

Lai & Robbins '85

P.R. Kumar '83

$H > 0$ fixed horizon, episodic setting, $0 \leq \delta < 1$

$$M = (S, A, P^*, r)$$

$\underbrace{P^*}_{\text{knows}} \quad \text{does not know}$

$$a \vee b = \max(a, b)$$

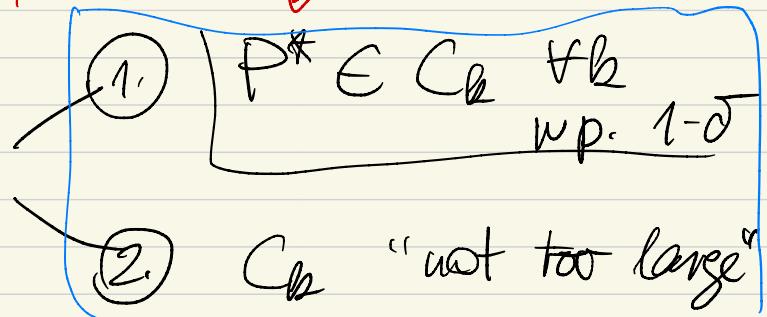
$$P_a^{(k)}(s, s') = \frac{N_k(s, a, s')}{1 \vee N_k(s, a)} = 0$$

$$C_k = \{ P = (P_{\alpha}(s)) \mid \forall s, a : \| P_{\alpha}^{(k)}(s) - P_{\alpha}(s) \|_1 \leq \beta \underbrace{\text{N}_{\alpha}(a)}_{=0} \}$$

$\beta/0 = 1$

$$\beta : \mathbb{N} \rightarrow (0, \infty)$$

Goal of choosing β :



$$\pi_k = \operatorname{argmax}_{\pi} V_{P_k}^{\pi}(S_0^{(k)})$$

$$\rightarrow \tilde{P}_k = \operatorname{argmax}_{P \in C_k} V_P^*(S_0^{(k)})$$

Follow π_k up to the end of the

episode.

Step 1: $C_k = ?$

Step 2: $R_{k2} \leq ?$

s, a fixed

$$n \in \mathbb{N} \quad \# \text{ visits to } s, a \\ P_{n,a}(s, s') = \frac{\sum_{v=1}^n \mathbb{I}(X_v = s')}{n}$$

$X_n \in S$: next state observed
upon visiting (s, a)
the n^{th} time

Markov property $\Rightarrow (X_n)_{n=1}^{\infty}$ i.i.d.

$$P_a^*(s)$$

$$\|P\|_1 = \sup_{\|x\|_\infty \leq 1} \langle p, x \rangle$$

$$\|P_{n,a}(s) - P_a^*(s)\|_1 =$$

$$= \max_{x \in \{-1\}^S} \langle P_{n,a}(s) - P_a^*(s), x \rangle$$

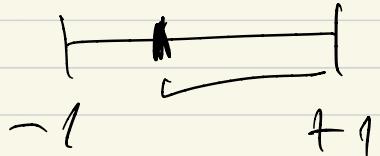
$$\text{Fix } x \in \{-1\}^S$$

$$\langle P_{\alpha,\alpha}(\zeta) - P_\alpha^*(\zeta), \chi \rangle =$$

$$\frac{1}{N} \sum_{n=1}^N \left[\sum_{\zeta' = 1}^S x_{\zeta'} (\mathbb{I}(X_n = \zeta) - P_\alpha^*(\zeta, \zeta')) \right]$$

Δ_n

$$= \frac{1}{n} \sum_{n=1}^N \Delta_n$$



$$\mathbb{E} \Delta_n = 0, \quad |\Delta_n| \leq \underline{2}$$

$(\Delta_n)_{n=1}^N$ i.i.d.

Hoeffding's \leq : $\forall \delta \in (0, 1)$

up $1 - \delta$

$$\frac{1}{n} \sum_{n=1}^N \Delta_n \leq 2 \sqrt{\frac{\log(1/\delta)}{2n}}$$

$$x \in \{-1\}^S$$

$$\{-1\}^S = 2^S$$

$$\mathcal{F} := \frac{\mathcal{O}}{2^S} \quad / \text{union bound} \Rightarrow$$

wp $\sim \delta$

$$\|P_{n,a}(s) - P_a^*(s)\|_1 \leq 2 \sqrt{\frac{s \log 2 + \log(1/\delta)}{2n}}$$

$\forall k: \sqrt{s, a}$

$$\|P_{n,a}^{(k)}(s) - P_a^*(s)\|_1 \leq \beta(N_k(s, a))$$

$$P_a^{(k)}(s) = P_{N_k(s, a), a}(s)$$

Take union bound over n, s, a

$$\sum_{n=1}^{\infty} \frac{\delta}{n(n+1)} = \delta$$

$$\sum \frac{1}{n^2} < \infty$$

$\rightarrow \text{wp} 1-\delta \quad \forall n \geq 1, \forall s, a$

$$\frac{1}{n} - \frac{1}{n+1}$$

$$\|P_{n,a}(s) - P_a^*(s)\|_1 \leq 2 \sqrt{\frac{s \log 2 + \log \frac{u(u+1)SA}{\delta}}{2n}}$$

$$\frac{\delta}{u(u+1)SA}$$

$$\beta(u) = 2 \sqrt{\frac{s \log 2 + \log \frac{u(u+1)SA}{\delta}}{2n}}$$

$\beta(u) \leq C_0/\sqrt{u}$, $C_0 = 2\sqrt{\frac{1}{2}(S \log 2 + \log \frac{4H(K+1)}{\delta})}$

Lemma 1: With β as above

$$\forall \delta \in (0, 1) \text{ wp } 1-\delta \quad P^* \in \bigcap_{k \geq 1} C_k$$

Fix δ , choose β as above.

Consider $\mathcal{E} = \{P^* \in \bigcap_{k \geq 1} C_k\}$

By Lemma 1: $P(\mathcal{E}) \geq 1-\delta$

In what follows assume that we are on \mathcal{E} .

$$R_K = \sum_{k=1}^K V_0^*(S^{(k)}) - V_K$$

$$V_K = \sum_{h=0}^{H-1} r_{A_h^{(k)}}(S_h^{(k)})$$

Fix $k \geq 1$.

$$V_0^*(S_0^{(k)}) - V_K$$

A

$$= V_{0, \tilde{P}^*}^*(S_0^{(k)}) - V_{0, \tilde{P}^{(k)}}(S_0^{(k)}) \quad \text{I.}$$

$$V_{0, \tilde{P}^{(k)}}^{\text{II}}(S_0^{(k)}) - V_{0, P^*}^{\text{II}}(S_0^{(k)}) \quad \text{II.}$$

$$V_{0, P^*}^{\text{III}}(S_0^{(k)}) - V_K \quad \text{III.}$$

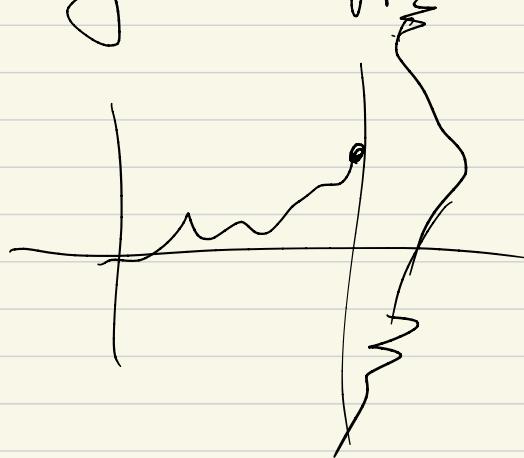
$$\text{I.} \leq 0 \quad P^* \in C_K$$

$$\text{II.} = V_{0, P^*}^{\text{II}}(S_0^{(k)}) - \sum_{h=0}^{H-1} r_{A_h^{(k)}}(S_h^{(k)})$$

$$\sum_{k=1}^K \xi_k \quad \left\{ \begin{array}{l} \xi_k \geq 0 \\ \sum_{k=1}^K \xi_k = 0 \\ \text{Rog}(S_K) \leq H \end{array} \right.$$

$$\sum_{k=1}^H \xi_k \leq H \sqrt{\frac{k}{2} \log(1/\delta)}$$

martingale difference seq (\approx i.i.d.)



+ which bond (~ 20)

Hoeffding - Azuma

$$\underline{D}_j = V_{0, P^*}^{\pi_k}(S_j^{(k)}) - V_{0, \widetilde{P}^{(k)}}^{\pi_k}(S_j^{(k)}) \quad (= \delta_j^{(k)})$$

$$V_{h, P}^{\pi} = r^{\pi} + M_{\pi} P V_{h+1, P}^{\pi}, \quad 0 \leq h \leq H-1$$

$$\boxed{V_{H, P}^{\pi} = 0}$$

$$\|P^* - P\|_{\infty}$$

$$V_{h, P^*}^{\pi} - V_{h, P}^{\pi} = M_{\pi} (P^* - P) V_{h+1, P^*}^{\pi}$$

$$\parallel \quad \|_\infty$$

$$M_{\bar{\pi}} P (V_{h+1,P^*}^{\bar{\pi}} - V_{h+1,P}^{\bar{\pi}})$$

$$(k) \quad V_{h,P^*}^{\bar{\pi}} - V_{h,P}^{\bar{\pi}} = M_{\bar{\pi}} P^* (V_{h+1,P^*}^{\bar{\pi}} - V_{h+1,P}^{\bar{\pi}}) \\ + M_{\bar{\pi}} (P^* - P) V_{h+1,P}^{\bar{\pi}}$$

$$V_{h,P^*}^{\bar{\pi}_k} (S_h^{(k)}) - V_{h,\tilde{P}^{(k)}}^{\bar{\pi}_k} (S_h^{(k)}) = \delta_h^{(k)} - \\ \underbrace{E[\delta_{h+1}^{(k)} | \mathcal{F}_{h,k}]}_{\mathcal{D}_{h+1}^{(k)}}$$

$$\delta_h^{(k)} = (M_{\bar{\pi}_k} P^* (V_{h+1,P^*}^{\bar{\pi}_k} - V_{h+1,\tilde{P}^{(k)}}^{\bar{\pi}_k})) (S_h^{(k)})$$

$$+ (M_{\bar{\pi}_k} (P^* - \tilde{P}^{(k)}) V_{h+1,\tilde{P}^{(k)}}^{\bar{\pi}_k}) (S_h^{(k)})$$

$$\leq \delta_{h+1}^{(k)} + (E[\delta_{h+1}^{(k)} | \mathcal{F}_{h,k}] - \delta_{h+1}^{(k)})$$

$$+ \left| P_{A_h^{(k)}}^* (S_h^{(k)}) - \tilde{P}_{A_h^{(k)}}^* (S_h^{(k)}) \right|_H$$

$$(M_{\pi}(P^* - P)v)(s) \quad \pi(s) = a$$

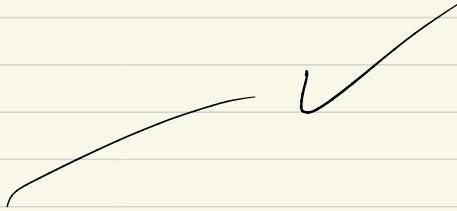
$$= \langle P_a^*(s) - P_a(s), v \rangle$$

$$\leq \|P_a^*(s) - P_a(s)\|_1 \|v\|_\infty$$

$$\leq \sigma_{h+1}^{(k)} + \gamma_{h+1}^{(k)} + 2\beta \left(N_h(S_h, A_h^{(k)}) \right)$$

↓ ↓

$$\sigma_H^{(k)} = 0$$



$$\sigma_0^{(k)} \leq \overbrace{\gamma_1^{(k)} + \dots + \gamma_{H-1}^{(k)}} +$$

$$+ H 2 \sum_{h=0}^{H-1} \beta \left(N_h(S_h, A_h^{(k)}) \right)$$

↓