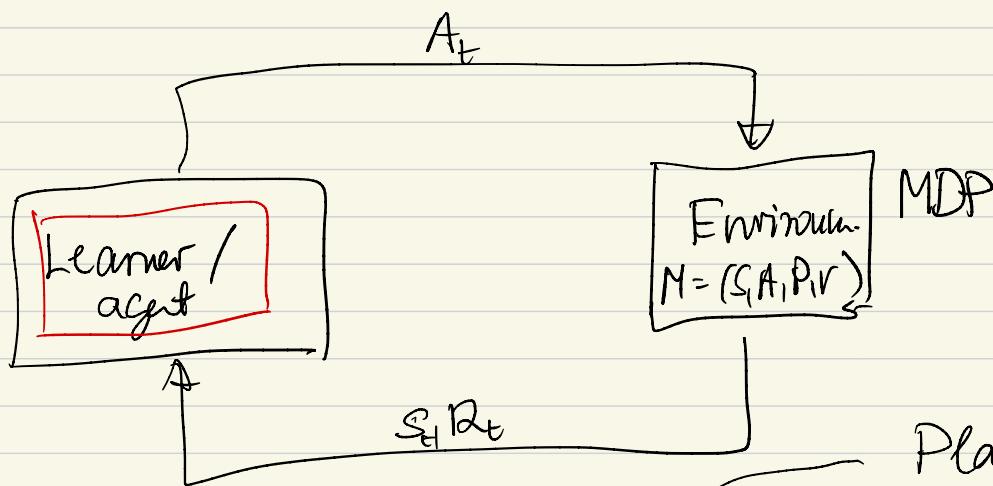
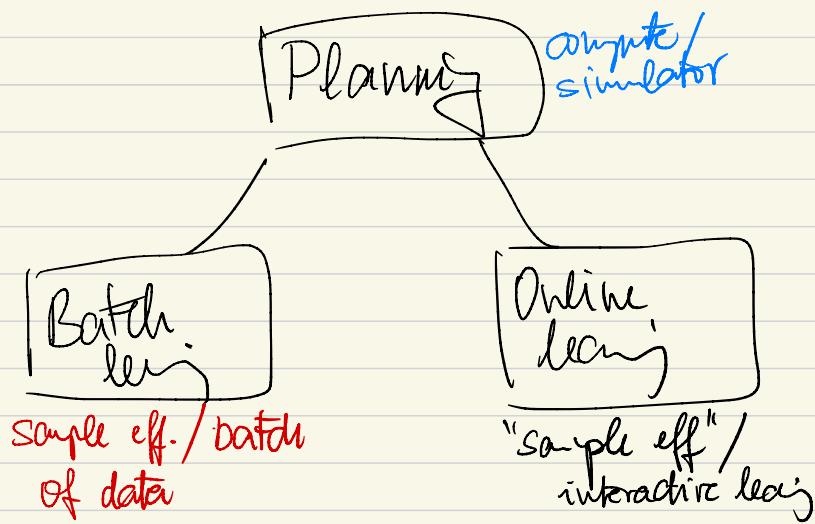


April 1

Online learning

/ Online RL



How good is a learner?

"Regret"

Cumulative regret

over time

undiscounted

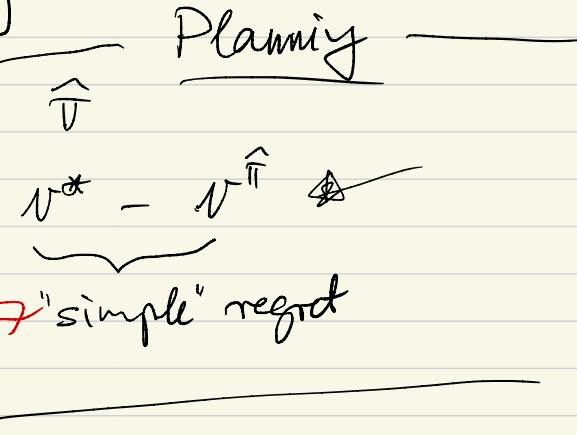
fixed-horizon

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

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

$$\underbrace{\frac{R_K}{K}}_{\text{regret}}$$

PAC \hookrightarrow Regret

MDP

$H = 1 \Leftrightarrow$ contextual bandits

$$|S|=1$$

finite-armed bandits

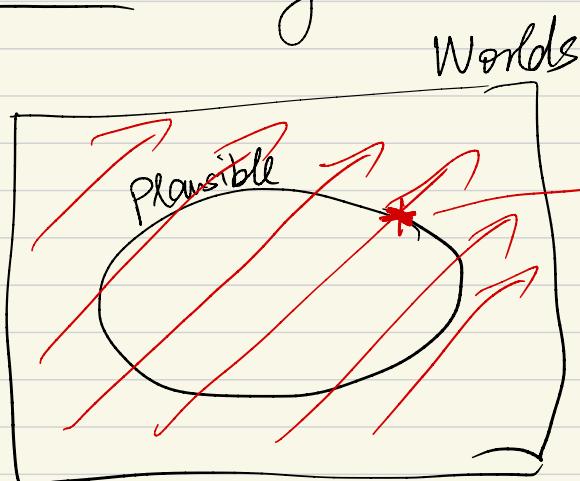
reward is unknown

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

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, $0 \leq \delta < 1$

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

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

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

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

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

$\beta / D = 1$

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

Goal of choosing β :

- 1. $P^* \in C_k \quad \forall k$
w.p. $1 - \delta$
- 2. C_k "not too large"

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

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

Follow π_k up to the end of the

episode.

Step 1: $C_k = ?$

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