

EEML

"Bits of turtles about RL"

1.

$$\mathcal{Q}(\sqrt{SA_n})$$

contextual bandits

2.

Optimism

3.

Function approx

4.

Batch RL

$$\mathcal{O}(A^{\min(SH)}/\epsilon^2)$$

Contextual bandit

Bandit = 1 state MDP

Contextual
bandits = S states, $S \sim D$ i.i.d.
 $A > 0$ # actions

Online

$1 \leq i \leq A$: actions

$\underline{\mu}_i(s, \cdot)$ reward distr. for action i .

$t = 1, 2, \dots$

I_t : choice of learner \leftarrow

$s_t \sim p(\cdot)$

$r_t \sim \mu_{I_t}(s_t, \cdot)$

} observed by learner

$$r_i(s) = \int r \mu_i(s, r) dr$$

$$\text{Regret}_n = \sum_{t=1}^n \max_i r_i(s_t) - R_t \quad R_t \in [0, 1]$$

$$\mathbb{E} [\text{Regret}_n] = R_n(t, \epsilon)$$

$$R_n^*(\epsilon) = \inf_{\mathcal{A}} \sup_{\epsilon \in \mathcal{C}} R_n(\mathcal{A}, \epsilon)$$

$$\mathcal{C} \quad \mu_i(s_i) \quad \in [0, 1] \quad \text{gaussian} \quad 1-\text{variance}$$

$$R_n^*(\epsilon) = \Theta(\sqrt{SA n})$$

$$\frac{R_n^*(\epsilon)}{n} = \Theta\left(\sqrt{\frac{SA}{n}}\right)$$

Part II

Explore - exploit dilemma

Optimism

Optimism in the face
of uncertainty

$$S = 1$$

r_1, \dots, r_A : means

$\hat{r}_1, \dots, \hat{r}_A$: empirical mean

$n_1, \dots, n_A > 0$: counts

UCB

$$\arg\max_{1 \leq i \leq A} \hat{r}_i + C \sqrt{\frac{\log(1/\delta)}{n_i}}$$

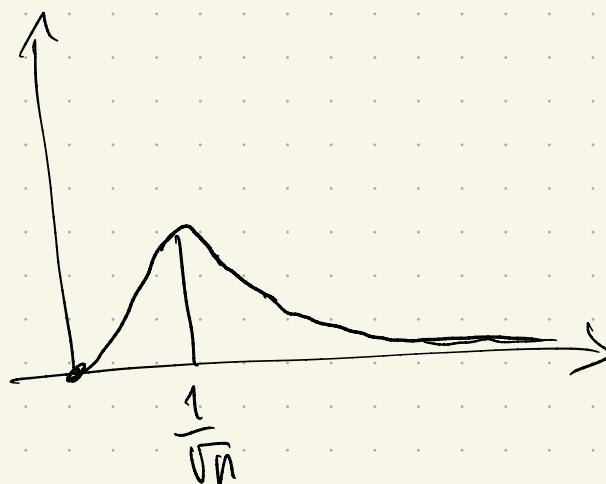
\tilde{r}_i

$$C > 0$$

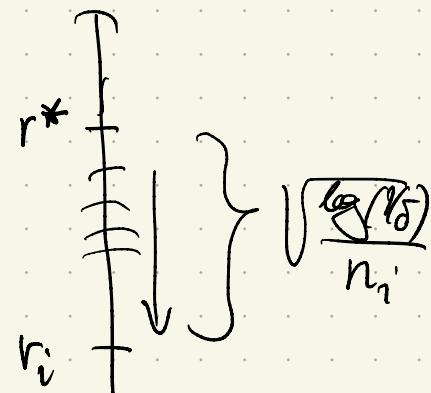
$$\delta = \frac{1}{n}$$

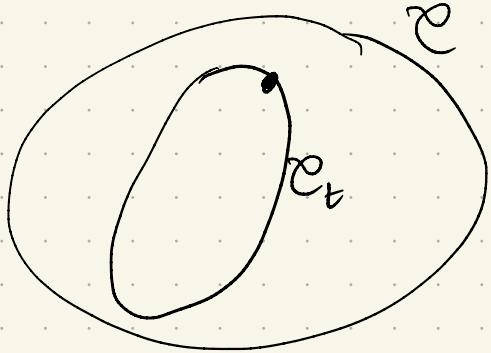
$$P(\tilde{r}_i < r_i) \approx \frac{1}{n}$$

Repet



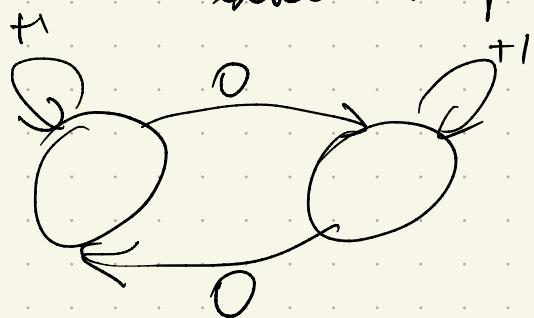
ϵ reward gap





controlled state

cozy policy has to
be used for an
extended period



Part 3

Function approximation

$$\mathcal{D}(\sqrt{SA_n})$$

Bellman

$$([S], [A], Q)$$

$$Q^*(s, a) = r(s, a) + \gamma$$

$$\sum_{s'} P(s, a, s') \max_{a'} Q^*(s')$$

$$\arg\max_a Q^*(s, a)$$

$$[S] = \{s_1, \dots, s\} \quad Q \approx Q^*$$

$$[A] = \{1, \dots, A\} \quad \arg\max_a Q(s, a)$$

$$q: [S] \times [A] \rightarrow \mathbb{R}^d \quad \exists \theta^* \in \mathbb{R}^d$$

$$Q^*(s, a) \approx q(s, a)^T \theta^*$$

$$\theta^* = ?$$

$$\hat{\theta} \approx \theta^*$$

Simulation / optimization / play

$$\arg\max_a q(s, a)^T \hat{\theta}$$

$$\text{Thm: } \forall \text{ planner } \mathcal{P} \quad \exists (\text{MDP}, q) \quad Q^* = q^T \theta^*$$

s.t. #queries used by \mathcal{P}

3

$$2 \mathcal{D}(\min(H, d))$$

$$H = \frac{1}{1-\gamma}$$

$$V^* = \varphi^T \theta^*$$

+ positive

$$O\left(\left(\frac{dH}{\delta}\right)^A\right)$$

5 Suboptimality of π induced by
planner

TD

$$A = 2$$

$$Q^\pi = \varphi^T \theta_\pi$$

poly time

policy iteration

policy eval

TD methods?

value iteration

$$T \mathcal{F} \subseteq \mathcal{F}$$

"closedness"

(Part IV.)

Batch RL

Off-policy
policy
opt.

Thm: Data generated by following some policy
query complexity = $\Theta(\underline{A}^{\min(H, S)} / \epsilon^2)$

