

$$\hat{P}_{ij}(a) = \frac{\sum_1^t \mathbb{1}_{\{x_t=i, A_t=a, x_{t+1}=j\}}}{\sum_1^t \mathbb{1}_{\{x_t=i, A_t=a\}}}$$

$\downarrow$   
 $T_t(i, a)$

$$M_t(i, a) = \left\{ P_{ij}(a) : \|P_{ij}(a) - \hat{P}_{ij}(a)\| \leq \frac{\sqrt{\log(4SAT_t(i, a))}}{\sqrt{\frac{2(1+T_t(i, a))}{\delta}}} \right\}$$

$\downarrow$   
 $M_t$ : Set of plausible MDPs.

In phase  $k$ , starting at  $T_k$ , play policy that is greedy w.r.t  $(P_k, h_k)$  where

$$P_k + h_k(i) = \min_a \left( c(i, a) + \sum_j \tilde{P}_{ij}(a) h(j) \right) \quad \text{for some } \tilde{P} \in M_{T_k} \quad (1)$$

$$P_k = \min_{\pi \in \Pi} \min_{M \in M_{T_k} \text{ stationary}} P^\pi(M)$$

Phase ends when  $T_t(x_t, A_t) = 2T_{T_k}(x_t, A_t)$ .

Extended MDP:

$$\tilde{S} = S$$

$$\tilde{A} = \{ (a, M) : a \in A, M \in M_{T_k} \}$$

$$\tilde{c}(i, (a, M)) = c(i, a)$$

$$\tilde{P}_{ij}(a, M) = P_{ij}^M(a)$$



↳ Solving above MDP, solves  $\textcircled{1}$ .  $\sum_{i=1}^n \frac{1}{i} = (n) \log n$

$M_\gamma$ : convex polyhedral set in higher dim.

Thm:

Under UCRL2,

$$E(R_n) \leq 1 + \frac{C D(M) \sqrt{2 A n \log n}}{1}$$

↓  
diameter of MDP.

univ. constant.  $\dots$   $D(M)$  directly to  $\sqrt{2 A n \log n}$

\* in case of BPI, if all mdp's in  $M$  have same action for each state, stop.

Sampling: not as trivial.

$\textcircled{1}$   $\left\{ \begin{array}{l} \text{if } M \neq \emptyset \text{ and } \dots \end{array} \right.$

$(M) \log \dots$

$(M, \gamma) \dots$

MDP:  $\dots$

$$2 = \sum$$

$$\{M \in M \mid A \in (M, A)\} = \tilde{A}$$

$$(M) \log = (M, A) \log$$

$$(M) \log = (M, A) \log$$

16/4/24

RL: Regret guarantees

Avg cost MDP:  $(S, A, P, c)$

↑ ↑  
finite

Under "suitable" conditions,  $\exists (P, h)$  satisfying

$$P + h(i) = \min_a (c(i, a) + \sum_j P_{ij}(a) h(j))$$

Policy greedy w.r.t  $(P, h)$  is optimal.

↳ choose  $a$  which  
minimizes RHS;  
for every  $i$ .

↳ max  $J$   
 $(P, h)$

$$c(i, a) +$$

$$s.t. \quad P + h(i) \leq \sum_j P_{ij}(a) h(j) \quad \forall i$$

$$h(0) = 0.$$

↳ reference state. (can be translated additively  
so anchor one state to  
zero).

↳ LP solution solves BE.

↳ RL:

$P$  or  $C$  or both unknown.

Regret  
minimization

BPI

(best policy identification).



↳ Regret minimisation:

$$\text{Minimize } \mathbb{E} \left( \sum_t C(X_t, A_t) \right) - n p^*$$

↓  
optimal  
time averaged  
cost  
(asymptotic)

↳ Turns out, Under optimal policy

$$\text{cost} = n p^* + O(1)$$

\* UCRL 2 :

(Upper Confidence Bound For RL)

Assume  $C$  is known.

Unknown:  $P$ .



confidence set: (plausible MDPs)  
 $\mathcal{M}_{\text{set}}$

$S, A, C$  fixed.

$P$  is the only unknown.

MDP  
with  
least  $p^*$ : optimal.

Play optimal action corresponding to best plausible MDP