

$$R_{\max} \quad M = (S, A, P, R, \gamma, d_0, H). \quad \frac{R_{\max}}{1-\gamma}$$

$$J(\pi^*) - J(\hat{\pi}) \leq \underline{\underline{\epsilon}} \cdot \underline{\underline{V_{\max}}}$$

How many ep do we need?

Algorithm:

$n(s,a)$
 $n(s,a,s')$ } counters.

$$\hat{P}(s'|s,a) = \frac{n(s,a,s')}{n(s,a)}$$

threshold m . $\hat{M}_k = (S, A, \hat{P}_k, R_k, \gamma)$

$$\hat{P}_k(s'|s,a) = \begin{cases} n(s,a,s')/n(s,a) & \text{if } n(s,a) = m. \\ \mathbb{I}[s'=s] & \text{o.w.} \end{cases}$$

$\Delta (s,a) \in K$
 $(s,a) \notin K$

$$R_k(s,a) = \begin{cases} R(s,a) & \text{if } n(s,a) = m. \\ R_{\max} & \text{o.w.} \end{cases}$$

"Optimism in face of uncertainty".

Next ep: explore w/ $\pi_{\hat{M}_k}^*$ in the next ep data. if any $(s_n, a_n) \notin K$.

$s_1, a_1, \boxed{s_2, a_2}, \dots, s_n, a_n$

inc $n(s_n, a_n)$
 $n(s_n, a_n, s_{n+1})$

"tape" / "cache".

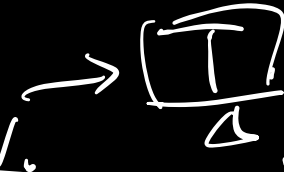
T_A .

A B A B B ...

fixed.

T_A
 \emptyset

A A A B A B A A ...



Y.V.

$$\forall s, a \in K \quad P(s'|s, a) = \begin{cases} P(s'|s, a) & \text{if } (s, a) \in K \\ \mathbb{I}[s'=s] & \text{otherwise} \end{cases}$$

Table 1: Relationship between M , M_K , and \widehat{M}_K .

	M	M_K	\widehat{M}_K
Known (K)	$= M$	$= M$	$\approx M$
Unknown	$= M$	self-loop	self-loop

$$M_K \quad P_K(s'|s, a) = \begin{cases} P(s'|s, a) & \text{if } (s, a) \in K \\ \mathbb{I}[s'=s] & \text{otherwise} \end{cases}$$

R_K

$$\text{w.p. } \geq 1 - \delta, \forall (s, a) \in K, \|\widehat{P}_K(\cdot | s, a) - P_K(\cdot | s, a)\|_1$$

$$\leq \underbrace{O\left(\sqrt{\frac{|S|}{m} \log \frac{|S \times A|}{\delta}}\right)}_{\epsilon_p}$$

$$\Rightarrow \forall \pi: S \rightarrow A.$$

$$\rightarrow \left\| V_{M_K}^\pi - V_{\widehat{M}_K}^\pi \right\|_\infty \leq \epsilon_p \cdot \frac{V_{\max}}{1 - \gamma}$$

$$\rightarrow \|V_{M_k}^* - \hat{V}_{M_k}^*\|_\infty \leq \epsilon_p \cdot \frac{V_{max}}{1-\gamma}$$

"Explore-or-Terminate"

Fact: optimism:

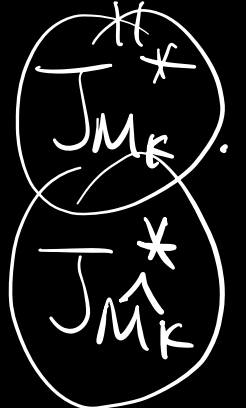
$$\forall \pi: S \rightarrow A.$$

$$\underline{J_M(\pi)} \leq \underline{J_{M_k}(\pi)}$$

$$\underline{\epsilon \cdot V_{max}} < \underline{J_M(\pi_M^*)} - \underline{J_M(\hat{\pi}_{M_k}^*)}$$

$$\leq J_{M_k}(\pi_M^*) - J_M(\hat{\pi}_{M_k}^*)$$

$$\leq \underline{J_{M_k}(\pi_{M_k}^*)} - \underline{J_M(\hat{\pi}_{M_k}^*)}$$



$$\leq \underline{J_{M_k}^*} - \underline{J_M(\hat{\pi}_{M_k}^*)} + \frac{\epsilon_p V_{max}}{1-\gamma}$$

$$\frac{\epsilon_p V_{max}}{1-\gamma}$$

$$= \underline{J_{M_k}^*}(\hat{\pi}) - \underline{J_M}(\hat{\pi}) + \frac{\epsilon_p V_{max}}{1-\gamma}$$

$$\frac{\epsilon_p V_{max}}{1-\gamma}$$

$$\frac{\epsilon V_{max}}{2} \leq$$

$$\underline{J_{M_k}^*}(\hat{\pi}) - \underline{J_M}(\hat{\pi}) + \frac{2\epsilon_p V_{max}}{1-\gamma}$$

$$\frac{2\epsilon_p V_{max}}{1-\gamma}$$

$$\frac{\epsilon V_{max}}{2}$$

$$\frac{\epsilon V_{max}}{2}$$

$$\frac{1}{1-\gamma} \mathbb{E}$$

$$[Q_{M_k}^{\hat{\pi}} - J_M Q_{M_k}^{\hat{\pi}}]$$

$$Q_{M_k}^{\hat{\pi}} - J_M Q_{M_k}^{\hat{\pi}}$$

$$= \sum_{(s,a) \in K} d_M^{\hat{\pi}}(s,a) \cdot (0) + \sum_{(s,a) \notin K} d_M^{\hat{\pi}}(s,a) \cdot (c \vee \max)$$

$$Q_{M,K}^{\hat{\pi}} - T_M Q_{M,K}^{\hat{\pi}} = T_{M,K} Q_{M,K}^{\hat{\pi}} - T_M Q_{M,K}^{\hat{\pi}}$$