

$$1. \sum_{s_t} |P_{\pi_\theta}(s_t) - P_{\pi^*}(s_t)| \quad (1)$$

$$\text{consider } P_{\pi_\theta}(s_t) = P_{\pi_\theta}(s_t | \text{no mistake}) \cdot P_{\pi_\theta}(\text{no mistake}) \\ + P_{\pi_\theta}(s_t | \text{mistake}) \cdot P_{\pi_\theta}(\text{mistake})$$

$$P_{\pi_\theta}(s_t | \text{no mistake}) = P_{\pi^*}(s_t), \quad P_{\pi_\theta}(\text{no mistake}) = 1 - P_{\pi_\theta}(\text{mistake})$$

$$\text{Consider } P_{\pi_\theta}(\text{mistake})$$

$$= P_{\pi_\theta} \cdot \left(\bigcup_{t' \in [1, t]} \{ \pi_\theta \text{ makes mistake at } t', \text{ but not before} \} \right)$$

$$\leq \sum_{t' \in [1, t]} P_{\pi_\theta} [\pi_\theta \text{ makes mistake at } t', \text{ but not before}]$$

Consider a Markov chain of three events

① π_θ makes mistake at t'

② s_t

③ π_θ makes no mistake before t'

$$\Rightarrow \textcircled{3} \rightarrow \textcircled{2} \rightarrow \textcircled{1}$$

$$\sum_{t' \in [1, t]} P_{\pi_\theta} [\pi_\theta \text{ makes mistake at } t', \text{ but not before}]$$

$$= \sum_{t' \in [1, t]} \sum_{s_t} P_{\pi_\theta} (\textcircled{1} \wedge \textcircled{2} \wedge \textcircled{3})$$

$$= \sum_{t' \in [1, t]} \sum_{s_t} P_{\pi_\theta} (\textcircled{1} | \textcircled{2}) P_{\pi_\theta} (\textcircled{2} | \textcircled{3}) \cdot P_{\pi_\theta} (\textcircled{3}) \quad (2)$$

$$\text{Note that } P_{\pi_\theta} (\textcircled{1} | \textcircled{2}) = \pi_\theta(a_{t'} \neq \pi_\theta(s_{t'}) | s_{t'})$$

$$P_{\pi_\theta} (\textcircled{2} | \textcircled{3}) = P_{\pi^*}(s_t)$$

$$P_{\pi_\theta} (\textcircled{3}) < 1$$

$$\therefore (2) \leq \sum_{t' \in [1, t]} E_{P_{\pi^*}(s_t)} \pi_\theta(a_{t'} \neq \pi_\theta(s_{t'}) | s_{t'})$$

$$\leq \sum_t E_{P_{\pi^*}(s_t)} \cdot \pi_\theta(a_t \neq \pi_\theta(s_t) | s_t)$$

$$\begin{aligned}
 (1) &= \sum_{s_t} P_{\pi_\theta}(\text{Mistake}) \cdot |P_{\pi_\theta}(s_t | \text{Mistake}) - P_{\pi^*}(s_t)| \\
 &\leq \underbrace{\sum_t E_{P_{\pi^*}(s_t)} \pi_\theta(a_t \neq \pi_\theta(s_t) | s_t)}_{\leq \varepsilon T} \cdot \underbrace{\sum_{s_t} |P_{\pi_\theta}(s_t | \text{Mistake}) - P_{\pi^*}(s_t)|}_{\leq 2} \\
 &\leq 2 \varepsilon T
 \end{aligned}$$

□

2.

$$\begin{aligned}
 (a) \quad J(\pi^*) - J(\pi_\theta) &= E_{P_{\pi^*}(s_T)} r(s_T) - E_{P_{\pi_\theta}(s_T)} r(s_T) \\
 &= \sum_{s_T} (P_{\pi^*}(s_T) - P_{\pi_\theta}(s_T)) \cdot r(s_T) \\
 &\leq r(s_T) \cdot \sum_{s_T} |P_{\pi^*}(s_T) - P_{\pi_\theta}(s_T)| \\
 &\leq R_{\max} \cdot 2T\varepsilon
 \end{aligned}$$

$$\therefore J(\pi^*) - J(\pi_\theta) = O(T\varepsilon) \quad \square$$

$$\begin{aligned}
 (b) \quad J(\pi^*) - J(\pi_\theta) &= \sum_{t=1}^T E_{P_{\pi^*}(s_t)} r(s_t) - \sum_{t=1}^T E_{P_{\pi_\theta}(s_t)} r(s_t) \\
 &= \sum_{t=1}^T \sum_{s_t} (P_{\pi^*}(s_t) - P_{\pi_\theta}(s_t)) \cdot r(s_t) \\
 &\leq T \cdot R_{\max} \cdot 2T\varepsilon
 \end{aligned}$$

$$\therefore J(\pi^*) - J(\pi_\theta) = O(T^2\varepsilon)$$