

I. Analysis

$$1. E_{\rho_{\pi^*(s)}} \pi_{\theta}(a \neq \pi^*(s) | s) \leq \epsilon$$

$$\Rightarrow \sum_s p_{\eta^*}(s) \pi_{\theta}(a \pm \pi^*(s) | s) \leq \epsilon$$

Now we consider that the event when the learned policy disagree with the expert policy at each time step is E

\Rightarrow The probability that at least one mistake occur over the horizon T is $\Pr[U_T E_T]$

Using the inequality: $\Pr[U_i E_i] \leq \sum_i \Pr[E_i]$
 $\Rightarrow \Pr[U_T E_T] \leq \sum_{E_i=1} \Pr[E_i]$

However: $\Pr[E] \leq \epsilon$
 $\Rightarrow \Pr[V_r E_T] \leq \sum_{t=1}^T \Pr[E_t] \leq \epsilon T$

$$\Rightarrow \sum_{s_t} |p_{\pi_\theta}(s_t) - p_{\pi^*}(s_t)| \leq 2\epsilon T$$

2.

$$\begin{aligned} \alpha / J(\pi^*) - J(\pi_\theta) &= \sum_{t=1}^T \left(\mathbb{E}_{p_{\pi^*}(s_t)} r(s_t) - \mathbb{E}_{p_{\pi_\theta}(s_t)} r(s_t) \right) \\ &= \sum_{t=1}^T \sum_{s_t} \left(p_{\pi^*}(s_t) r(s_t) - p_{\pi_\theta}(s_t) r(s_t) \right) \\ &= \sum_{t=1}^T \sum_{s_t} r(s_t) (p_{\pi^*}(s_t) - p_{\pi_\theta}(s_t)) \\ &\leq \sum_{t=1}^T |r(s_t)| |p_{\pi^*}(s_t) - p_{\pi_\theta}(s_t)| \quad \text{Proving in lecture: } |p_{\pi^*}(s_t) - p_{\pi_\theta}(s_t)| \leq 2\epsilon t \\ &\leq \sum_{t=1}^T |r(s_t)| 2\epsilon t \end{aligned}$$

With reward only exist in the last state

$$\Rightarrow J(n^*) - J(\pi_\theta) \leq \sum_{t=1}^T |r(s_t)| 2\epsilon t$$
$$= 2\epsilon T$$
$$= O(\epsilon T) \quad \text{"Proved!!!"} \quad \square$$

b/ Similar to a/, We have: $J(\pi^*) - J(\pi_\theta) \leq \sum_{t=1}^T |r(s_t)| 2\epsilon T$. Since this question ask for an arbitrary reward

$$\leq \sum_{t=1}^T R_{\max} 2\epsilon T = 2\epsilon T^2 R_{\max} = O(\epsilon T^2) \quad \text{"Proved !!!"}$$