

## 1.4 & 1.5 Sampling Models & The Performance Difference Lemma

## 1. The episodic setting.

Learners acts for finite steps, they start from a fixed starting state  $s_0 \sim \mu$ , observe the trajectory and then reset to  $s_0 \sim \mu$ .

(i) Finite Horizon MDPs:

each episode lasts for  $H$  steps.

(ii) Infinite Horizon MDPs:

★1: agents can terminate episodes after fixed steps;

★2: each step has a probability of  $1-\gamma$  to terminate.  
this leads to an unbiased estimate of  $V^\pi$ .

Interests:

(i) number of episodes to find a near optimal policy

(ii) regret guarantee.

(iii) the strategy for the agents' exploration.

## 2. The generative model setting

Input a state-action pair  $(s, a)$  ;

Return a sample  $s' \sim P(\cdot | s, a)$  and  $r(s, a)$

## 3. The offline RL setting.

Agents has access to an offline dataset generated under certain policy.

under certain policy.

we assume the dataset is of the form  $\{(s, a, s', r)\}$   
 $s' \sim P(\cdot | s, a)$ ,  $r \sim r(s, a)$ .  $(s, a) \sim \Delta(S \times A)$  i.i.d.

### 1.5 The performance difference lemma.

## [Notations]

$$\textcircled{1} \quad V^\pi(\mu) := \mathbb{E}_{s \sim \mu}[V^\pi(s)]$$

② advantage :  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$   
 $A^*(s, a) := A^{\pi^*}(s, a) \leq 0$

[Def] visitation measure over states :

$$d_{S_0}^\pi(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P_r^\pi(s_t=s|s_0)$$

$$d_{\mu}^{\pi}(s) = \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^{\pi}(s)]$$

note:  $\sum_s d_\mu^\pi(s) = E_{s_0 \sim \mu} [\sum_s d_{s_0}^\pi(s)] = 1$ .  $d_\mu^\pi$  is a distribution.

Lemma 1.16 [The performance lemma]

For all policies  $\pi, \pi'$  and distributions  $\mu$  over  $S$ :

$$V^{\pi}(s) - V^{\pi'}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_{\mu}^{\pi}} \mathbb{E}_{a' \sim \pi(s')} [A^{\pi'}(s', a')]$$

Pf: Let  $P_r^\pi(\tau | s_0 = s)$  denote the probability of observing a trajectory  $\tau$  when starting in state  $s$  and following  $\pi$ .

① first we show that:

$$E_{T \sim p_T^\pi(s_t)} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \frac{1}{1-\gamma} E_{s \sim d_{\zeta}^\pi} E_{a \sim \pi(s)} [f(s, a)]$$

18

$$\begin{aligned}
\textcircled{2} \quad V''(s) - V''(s) &= E_{T \sim P_r^\pi(T|S_0=s)} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] - V''(s) \\
&= E_{T \sim P_r^\pi(T|S_0=s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + V^{\pi'}(s_t) - V^{\pi'}(s_t)) \right] - V^{\pi'}(s) \\
&= E_{T \sim P_r^\pi(T|S_0=s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \gamma V^{\pi'}(s_{t+1}) - V^{\pi'}(s_t)) \right] \\
&= E_{T \sim P_r^\pi(T|S_0=s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \gamma E[V^{\pi'}(s_{t+1})|s_t, a_t] - V^{\pi'}(s_t)) \right] \\
&= E_{T \sim P_r^\pi(T|S_0=s)} \left[ \sum_{t=0}^{\infty} \gamma^t (Q^{\pi'}(s_t, a_t) - V^{\pi'}(s_t)) \right] \\
&= E_{T \sim P_r^\pi(T|S_0=s)} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi'}(s_t, a_t) \right] \\
&= \frac{1}{1-\gamma} E_{s' \sim d_s^\pi} E_{a \sim \pi(\cdot|s)} [A^{\pi'}(s', a)] \quad \square
\end{aligned}$$