## 2.2 Sublinear Sample Complexity

Previously, we are able to estimate value of any policy.
Now we focus on the estimate of $\hat{Q}^*$.

#samples from generative model $= |S||A| N$

---

**Lemma 2.5** [component-wise bounds]
$$Q^* - \hat{Q}^* \leq \gamma(I - \gamma \hat{P}^{\pi^*})^{-1}(P - \hat{P})V^* \quad (1)$$
$$Q^* - \hat{Q}^* \geq \gamma(I - \gamma \hat{P}^{\hat{\pi}})^{-1}(P - \hat{P})V^* \quad (2)$$

---

Pf: (1): $\pi^*$ is optimal for $M$, $\hat{\pi}$ is optimal for $\hat{M}$
$$Q^* - \hat{Q}^* = Q^{\pi^*} - \hat{Q}^{\hat{\pi}}$$
$$\leq Q^{\pi^*} - \hat{Q}^{\pi^*}$$
$$(\text{Lemma 2.2}) = \gamma(I - \gamma \hat{P}^{\pi^*})^{-1}(P - \hat{P})V^{\pi^*}$$

(2)
$$Q^* - \hat{Q}^* = Q^{\pi^*} - \hat{Q}^{\hat{\pi}}$$
$$= \gamma(I - \gamma \hat{P}^{\hat{\pi}})^{-1}(P^{\pi^*} - \hat{P}^{\hat{\pi}})Q^*$$
$$\geq \gamma(I - \gamma \hat{P}^{\hat{\pi}})^{-1}(P^{\pi^*} - \hat{P}^{\pi^*})Q^*$$
$$= \gamma(I - \gamma \hat{P}^{\hat{\pi}})^{-1}(P - \hat{P})V^*$$
where $\hat{P}^{\hat{\pi}}Q^* \leq \hat{P}^{\pi^*}Q^*$ is because $\pi^*(s) = \arg\max_{a \in A} Q^*(s,a)$

---

**Proposition 2.4** [Crude Value Bounds]
Let $\delta \geq 0$, with prob. $\geq 1 - \delta$,
$$\|Q^* - \hat{Q}^*\|_\infty \leq \Delta_{\delta, N}$$
$$\|Q^* - \hat{Q}^{\pi^*}\|_\infty \leq \Delta_{\delta, N},$$
where

$$\Delta_{\delta,N} := \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2\log(2|S||A|/\delta)}{N}}$$

Pf: Due to Lemma 2.2 and Lemma 2.3

$$\|Q^* - \hat{Q}^{\pi^*}\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P-\hat{P})V^*\|_\infty \qquad (3)$$

Due to Lemma 2.5 and Lemma 2.3

$$\|Q^* - \hat{Q}^*\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P-\hat{P})V^*\|_\infty \qquad (4)$$

By applying Hoeffding's ineq.

$$\|(P-\hat{P})V^*\|_\infty = \max_{s,a} \left| E_{s'\sim P(s,a)}[V^*(s')] - E_{s'\sim \hat{P}(s,a)}[V^*(s')] \right|$$

$$\leq \frac{1}{1-\gamma} \sqrt{\frac{2\log(2|S||A|/\delta)}{N}} \quad ,$$

which holds with prob. $\geq 1-\delta$. $\qquad \square$

Addition :

Hoeffding's ineq. : $P(|S_n - E[S_n]| \geq t) \leq 2\exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i-a_i)^2}\right)$

where $a_i \leq \mathcal{I}_i \leq b_i$, $S_n = \sum_{i=1}^n \mathcal{I}_i$.

Since $E_{s'\sim \hat{P}(s,a)}[V^*(s')] = \sum_{i=1}^N V^*(s_i)/N$

and $0 \leq V^*(s_i)/N \leq \frac{1}{(1-\gamma)N}$, we have

$$P\left(|E_{s'\sim P(s,a)}[V^*(s^*)] - E_{s'\sim \hat{P}(s,a)}[V^*(s')]| \geq \frac{1}{1-\gamma}\sqrt{\frac{2\log(2|S||A|/\delta)}{N}}\right)$$

$$\leq 2\exp\left\{\frac{-2 \cdot \frac{1}{(1-\gamma)^2} \cdot \frac{2\log(2|S||A|/\delta)}{N}}{\frac{1}{(1-\gamma)^2} N^2}\right\}$$

$$= \exp\{-4N\} \frac{1}{|S||A|} \delta$$

$$\leq \delta$$

Actually $\frac{\gamma}{(1-\gamma)^2} \cdot \frac{\log(1/\delta)}{2N}$ for the estimate is enough.

Q: 2.1 seems to be sublinear as well.