

## 2.1 a naive model-based approach

Central question:

Do we require an accuracy model of the world to find a near optimal policy?

A naive model to learn  $P$ : after sampling  $N$  times, let  $\hat{P}(s'|s, a) = \frac{\text{count}(s', s, a)}{N}$

here we view  $\hat{P}$  as a matrix of size  $|S||A| \times |S|$

Expectation:  $O(|S|^2|A|)$  observation is enough for an accurate model.

Proposition 2.1 Assume  $\epsilon \in (0, \frac{1}{1-\gamma})$ ,  $\exists c > 0$  s.t.

# samples from generative model

$$= |S||A|N \geq \frac{4c^2}{(1-\gamma)^4} \frac{|S|^2|A|\log(|S||A|/\delta)}{\epsilon^2}$$

where  $(s, a)$  is sampled uniformly. and with prob.  $> 1-\delta$  we have

① (Model accuracy)

$$\max_{s, a} \|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \leq (1-\gamma)^2 \epsilon$$

② (Uniform value accuracy)

$$\|Q^\pi - \hat{Q}^\pi\|_\infty \leq \frac{\epsilon}{2} \quad \text{for all } \pi$$

③ (Near optimal planning) Suppose  $\hat{\pi}$  is optimal w.r.t.  $\hat{M}$

$$\|\hat{Q}^* - Q^*\|_\infty \leq \frac{\epsilon}{2}, \quad \|Q^{\hat{\pi}} - Q^*\|_\infty \leq \epsilon$$

To show this, we need following lemmas.

Lemma 2.2 [Simulation lemma] For all  $\pi$  :

$$Q^\pi - \hat{Q}^\pi = \gamma (I - \gamma \hat{P}^\pi)^{-1} (P - \hat{P}) V^\pi.$$

$$\begin{aligned} \text{Pf: } Q^\pi - \hat{Q}^\pi &= Q^\pi - (I - \gamma \hat{P}^\pi)^{-1} r \\ &= (I - \gamma \hat{P}^\pi)^{-1} ((I - \gamma \hat{P}^\pi) Q^\pi - r) \\ &= (I - \gamma \hat{P}^\pi)^{-1} ((I - \gamma \hat{P}^\pi) - (I - \gamma P^\pi)) Q^\pi \\ &= \gamma (I - \gamma \hat{P}^\pi)^{-1} (P^\pi - \hat{P}^\pi) Q^\pi \\ &= \gamma (I - \gamma \hat{P}^\pi)^{-1} (P - \hat{P}) V^\pi \quad \square \end{aligned}$$

Lemma 2.3 For any policy  $\pi$ , MDP  $M$  and  $v \in \mathbb{R}^{|S||A|}$

$$\|(I - \gamma P^\pi)^{-1} v\|_\infty \leq \frac{1}{1-\gamma} \|v\|_\infty$$

$$\text{Pf: } v = (I - \gamma P^\pi)(I - \gamma P^\pi)^{-1} v =: (I - \gamma P^\pi) w$$

$$\begin{aligned} \Rightarrow \|v\|_\infty &= \|(I - \gamma P^\pi) w\|_\infty \\ &\leq \|w\|_\infty - \gamma \|P^\pi w\|_\infty \\ &\leq \|w\|_\infty - \gamma \|w\|_\infty \\ &= (1-\gamma) \|w\|_\infty \end{aligned}$$

$$\text{i.e. } \|(I - \gamma P^\pi)^{-1} v\| \leq \frac{1}{1-\gamma} \|v\|_\infty \quad \square$$

Lemma A.8 [Concentration for discrete distributions]

Let  $z$  be r.v. of  $\{1, \dots, d\}$ , distributed according to  $q$ , where  $\bar{q} = [P_r(z=j)]_{j=1}^d$ . Assume we have  $N$  i.i.d. samples and that our empirical estimate is  $[\hat{q}]_j = \frac{\sum_{i=1}^N \mathbb{1}_{\{z_i=j\}}}{N}$ ,

we have  $\forall \epsilon > 0$ :

$$P_r(\|\hat{q} - q\|_2 \geq \frac{1}{\sqrt{N}} + \epsilon) \leq e^{-N\epsilon^2},$$

which implies:

$$P_r(\|\hat{q} - \bar{q}\|_1) \geq \sqrt{d} \left( \frac{1}{\sqrt{N}} + \varepsilon \right) \leq e^{-N\varepsilon^2}.$$

this proof is ignored

Pf of Proposition 2.1 :

with  $\ell_1$  norm in lemma A.8, for fixed  $s, a$ , with prob.  $\geq 1 - \delta$ , we have

$$\|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \leq c \sqrt{\frac{|S| \log(C1/\delta)}{N}} \quad (*)$$

where  $N$  is the number of samples used to estimate  $\hat{P}(\cdot|s, a)$ . Let  $|S|A|\delta = e^{-N\varepsilon^2} \Rightarrow \varepsilon = \sqrt{\frac{\log(C|S|A|/\delta)}{N}}$ ,  $d = |S|$  and let  $c$  satisfy  $c \sqrt{\frac{|S| \log(C|S|A|/\delta)}{N}} \geq \sqrt{|S|} \left( \frac{1}{\sqrt{N}} + \varepsilon \right)$

$$\textcircled{1} \|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \leq (1-\gamma)^2 \varepsilon$$

since  $N \geq \frac{4c^2}{(1-\gamma)^4} \frac{|S| \log(C|S|A|/\delta)}{\varepsilon^2}$ , by (\*) we have

$$\|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \leq (1-\gamma)^2 \varepsilon / 2 \quad \text{with prob. } \geq 1 - \delta$$

$$\textcircled{2} \|Q^\pi - \hat{Q}^\pi\|_\infty \leq \frac{\varepsilon}{2}$$

By Lemma 2.2:

$$\|Q^\pi - \hat{Q}^\pi\|_\infty = \|\gamma(I - \gamma P^\pi)^{-1}(P - \hat{P})V^\pi\|_\infty$$

$$\text{Lemma 2.3} \leq \frac{\gamma}{1-\gamma} \|(P - \hat{P})V^\pi\|_\infty$$

$$\begin{aligned} \text{H\"older ineq.} &\leq \frac{\gamma}{1-\gamma} \left( \max_{s, a} \|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \right) \|V^\pi\|_\infty \\ &\leq \frac{\gamma}{(1-\gamma)^2} \max_{s, a} \|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \\ &\leq \gamma \varepsilon / 2 \leq \varepsilon / 2 \end{aligned}$$

$$\textcircled{3} \|\hat{Q}^* - Q^*\|_\infty \leq \frac{\varepsilon}{2}, \quad \|Q^\pi - Q^{\pi^*}\|_\infty \leq \varepsilon$$

observe that  $|\sup_x f(x) - \sup_x g(x)| \leq \sup_x |f(x) - g(x)|$

$$\begin{aligned} \Rightarrow |\hat{Q}^*(s, a) - Q^*(s, a)| &= \left| \sup_{\pi} \hat{Q}^\pi(s, a) - \sup_{\pi} Q^\pi(s, a) \right| \\ &\leq \sup_{\pi} |\hat{Q}^\pi(s, a) - Q^\pi(s, a)| \end{aligned}$$

$$\begin{aligned}
\|Q^{\hat{\pi}} - Q^{\pi^*}\|_{\infty} &\leq \|Q^{\hat{\pi}} - \hat{Q}^*\|_{\infty} + \|\hat{Q}^* - Q^{\pi^*}\|_{\infty} \\
&= \|Q^{\hat{\pi}} - \hat{Q}^{\hat{\pi}}\|_{\infty} + \|\hat{Q}^* - Q^*\|_{\infty} \\
&\leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \quad \square
\end{aligned}$$