

Lecture 10: Model Based RL

MDP model  $\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, P, r, \gamma \}$  infinite horizon tabular  
 states actions transition dynamics reward cost discount  
 or  $\mathcal{M} = \{ \mathbb{R}^{n_s}, \mathbb{R}^{n_a}, f, c, H, \mu_0 \}$  finite horizon continuous.  
 initial distribution

But now transitions/dynamics are unknown!

### 1) MBRL Algorithm with query model

The query model (also called generative model):

For any  $s, a$  we can query the transition/dynamics model to sample the next state.

$$s' \sim P(s, a) \quad (\text{equivalently, } s' \sim f(s, a, w) \text{ s.t. } w \sim \mathcal{D})$$

Black-box sampling access.

Applicable to games + physics simulators.

Also simple, so it is a good starting point to understand sample complexity: How many samples are required for good performance?

### Alg: MBRL with query model

1) For  $i = 1, \dots, N$ :

Sample  $s'_i \sim P(s_i, a_i)$  and record  $(s'_i, s_i, a_i)$

2) Fit transition model  $\hat{P}$  from data  $\{(s'_i, s_i, a_i)\}_{i=1}^N$

3) Design  $\hat{\pi}$  using  $\hat{P}$

Today we will investigate the sample complexity of this method in two specific settings: tabular & LQR.

## 2) Tabular Setting

Specializing the algorithm to this setting:

1) sample all  $(s, a)$  evenly:  $\frac{N}{SA}$  times each

2) Fit transition model by counting

$$\hat{P}(s'|s, a) = \frac{\sum_{i=1}^N \mathbb{1}\{s_i = s \& a_i = a\} \mathbb{1}\{s'_i = s'\}}{\sum_{i=1}^N \mathbb{1}\{s_i = s \& a_i = a\}}$$

3) Design  $\hat{\pi}$  with Policy Iteration:  $\hat{\pi} = PI(\hat{P}, r)$

Recall:  $PI(P, r)$

Initialize  $\pi^0$

For  $t=1, \dots, T$ :

$Q^{\pi^t} = \text{Policy Eval}(\pi^t; P, r)$

$\pi^t(s) = \arg\max_a Q^{\pi^t}(s, a) \quad \forall s$

$$\begin{cases} V^{\pi} = (I - \gamma P)^{-1} R \\ Q^{\pi^t}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P} [V^{\pi}(s')] \end{cases}$$

Goal: Compare performance of  $\pi_*$  vs.  $\hat{\pi}$

strategy: i) compare  $P$  vs.  $\hat{P}$

ii) Translate  $P$  vs.  $\hat{P}$  into difference between value functions

iii) Translate difference in value functions to  $PI$

i)  $P$  vs.  $\hat{P}$ : similar to last lectures discussion

Lemma: with probability  $1-\delta$ , for all  $s, a$

$$\sum_{s' \in S} |\hat{P}(s'|s, a) - P(s'|s, a)| \leq \sqrt{\frac{S^2 A \log(2SA/\delta)}{N}}$$

$\|\hat{P}(\cdot|s, a) - P(\cdot|s, a)\|_1$

Proof is out of scope

## 11) Value Functions: effect of model error

Given a policy  $\pi$ , what is the difference between the value function defined by  $P$  compared to the value function defined by  $\hat{P}$ ?

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, s_{t+1} \sim P(s_t, a_t), a_t = \pi(s_t) \right]$$

$$\hat{V}^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, s_{t+1} \sim \hat{P}(s_t, a_t), a_t = \pi(s_t) \right]$$

Recall: Discounted state-action distribution

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P_t^\pi(s, a; s_0)$$

↑  
probability of visiting  $s, a$  at step  $t$   
starting at initial state  $s_0$

Simulation Lemma:

$$\hat{V}^\pi(s_0) - V^\pi(s_0) \leq \frac{1}{(1-\gamma)^2} \mathbb{E}_{s, a \sim d_{s_0}^\pi} \left[ \|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \right]$$

↑  
distribution under true  $P$

↑  
disagreement  $\hat{P}$  vs.  $P$

Proof: First, we claim that

$$\hat{V}^\pi(s_0) - V^\pi(s_0) = \gamma \mathbb{E}_{a_0 \sim \pi(s_0)} \left[ \mathbb{E}_{s_1 \sim \hat{P}(s_0, a_0)} [\hat{V}^\pi(s_1)] - \mathbb{E}_{s_1 \sim P(s_0, a_0)} [\hat{V}^\pi(s_1)] \right]$$

$$+ \gamma \mathbb{E}_{a_0 \sim \pi(s_0)} \left[ \hat{V}^\pi(s_1) - V^\pi(s_1) \right]$$

$s_1 \sim P(s_0, a_0)$

By iterating this expression  $K$  times,

$$\hat{V}^\pi(s_0) - V^\pi(s_0) = \sum_{t=1}^K \gamma^t \mathbb{E}_{s_{t-1}, a_{t-1}} \left[ \mathbb{E}_{s_t \sim \hat{P}(s_{t-1}, a_{t-1})} [\hat{V}^\pi(s_t)] - \mathbb{E}_{s_t \sim P(s_{t-1}, a_{t-1})} [\hat{V}^\pi(s_t)] \right]$$

$$+ \gamma^K \mathbb{E}_{s_K \sim P(s_{K-1}, a_{K-1})} [\hat{V}^\pi(s_K) - V^\pi(s_K)]$$

letting  $K \rightarrow \infty$ ,

$$\hat{V}^\pi(s_0) - V^\pi(s_0) = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[ \mathbb{E}_{s' \sim \hat{P}(s,a)} [\hat{V}^\pi(s')] - \mathbb{E}_{s' \sim P(s,a)} [V^\pi(s')] \right]$$

$$\mathbb{E}_{s' \sim \hat{P}} (\hat{V}^\pi(s')) - \mathbb{E}_{s' \sim P} (V^\pi(s')) = \sum_{s' \in \mathcal{S}} (\hat{P}(s'|s,a) - P(s'|s,a)) \hat{V}^\pi(s')$$

since  $r(s,a) \leq 1$ ,  $\hat{V}^\pi(s') \leq \frac{1}{1-\gamma}$

$$\leq \sum_{s' \in \mathcal{S}} |\hat{P}(s'|s,a) - P(s'|s,a)| \frac{1}{1-\gamma}$$

$$= \frac{1}{1-\gamma} \|\hat{P}(\cdot|s,a) - P(\cdot|s,a)\|_1$$

Then all that's left is to prove the initial claim.

$$\hat{V}^\pi(s_0) - V^\pi(s_0) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, \hat{P} \right] - \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, \hat{P} \right]$$

( $t=0$  term is equal)

$$= \gamma \mathbb{E}_{a_0 = \pi(s_0)} \left[ \mathbb{E}_{s_1 \sim \hat{P}(s_0, a_0)} [\hat{V}^\pi(s_1)] - \mathbb{E}_{s_1 \sim P(s_0, a_0)} [V^\pi(s_1)] \right]$$

$$= \gamma \mathbb{E}_{a_0 = \pi(s_0)} \left[ \underbrace{\mathbb{E}_{s_1 \sim \hat{P}} [\hat{V}^\pi(s_1)] - \mathbb{E}_{s_1 \sim P} [\hat{V}^\pi(s_1)]}_{\text{}} + \underbrace{\mathbb{E}_{s_1 \sim P} [\hat{V}^\pi(s_1)] - \mathbb{E}_{s_1 \sim P} [V^\pi(s_1)]}_{\text{}} \right] \quad \checkmark \quad \square$$

### III) Policy Iteration

Let  $\hat{\pi}^* = \text{PI}(\hat{P}, r) \leftarrow$  ignore iteration approximation  
for now - assume  $T > SA$  (HW1)

comparing to true optimal value.

$$V^*(s_0) - V^{\hat{\pi}^*}(s_0) \leq V^*(s_0) - \underbrace{\hat{V}^{\pi^*}(s_0)}_{\hat{\pi}^* \text{ is optimal on } \hat{P} \text{ so } \hat{V}^{\hat{\pi}^*}(s) \geq \hat{V}^{\pi}(s_0) \forall \pi.} + \hat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0)$$

(simulation lemma 2x)

$$\leq \frac{1}{(1-\gamma)^2} \left( \mathbb{E}_{s_1 \sim d_{s_0}^{\hat{\pi}^*}} \|\hat{P}(\cdot|s_1, a) - P(\cdot|s_1, a)\|_1 + \mathbb{E}_{s_1 \sim d_0^{\hat{\pi}^*}} \|\hat{P}(\cdot|s_1, a) - P(\cdot|s_1, a)\|_1 \right)$$

(model error bound)

$$\leq \frac{2}{(1-\gamma)^2} \sqrt{\frac{S \log(2SA/\delta)}{N}} \quad \text{w.p. } 1-\delta$$

Theorem: (sample complexity)

For  $0 \leq \delta \leq 1$ ,  $0 \leq \varepsilon \leq \frac{1}{1-\gamma}$ , let  $N = \frac{4S^2 A \log(\frac{2SA}{\delta})}{\varepsilon^2 (1-\gamma)^4}$

Then with probability at least  $1-\delta$ ,

$$V^*(s_0) - V^{\hat{\pi}^*}(s_0) \leq \varepsilon.$$

### 3) LQR

MBRL in this setting:

1) generate iid. samples  $s_i \sim \mathcal{N}(0, \sigma^2)$ ,  $a_i \sim \mathcal{N}(0, \sigma^2)$

2) estimate parameters by least squares

$$(\hat{A}, \hat{B}) = \operatorname{argmin} \sum_{i=1}^N (s'_i - A s_i - B a_i)^2$$

3) compute  $K_* = \text{LQR}(\hat{A}, \hat{B}, Q, R)$

We won't derive results in detail for this setting. But at a high level,

i) parameter estimation

$$\left\| \begin{bmatrix} \hat{A} - A \\ \hat{B} - B \end{bmatrix} \right\|_2 \lesssim \sqrt{\frac{(n_s + n_a) \log(1/\delta)}{N}}$$

matrix norm  $\nearrow$

ii) Difference in value ( $V_t^*(s) = s^T P_t s + p_t$ )

$$\|P_t - \hat{P}_t\|_2 \lesssim \left\| \begin{bmatrix} \hat{A} - A \\ \hat{B} - B \end{bmatrix} \right\|_2$$

iii) Difference in performance

$$\hat{V}_0^*(s_0) - V_0^*(s_0) \lesssim \|P_t - \hat{P}_t\|_2 \lesssim \sqrt{\frac{(n_s + n_a) \log(1/\delta)}{N}}$$

Sample complexity:  $\epsilon$ -optimal policy

$$\text{after } N \gtrsim \frac{(n_s + n_a)^2}{\epsilon^2} \text{ samples}$$