

# Lecture 18 Continuous State MDP & Model Simulation

1. Notation for continuous state MDP.

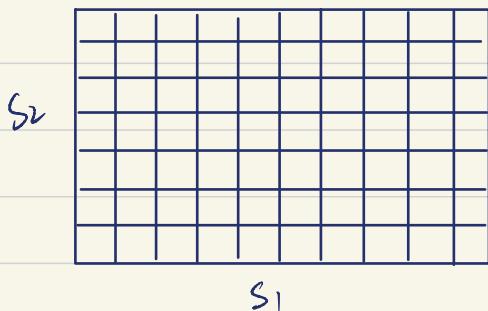
e.g. car :  $S(x, y, \theta, \dot{x}, \dot{y}, \dot{\theta})$

direction velocity on each direction

helicopter :  $S(x, y, z, \phi, \theta, \psi, \dot{x}, \dot{y}, \dot{z}, \dot{\phi}, \dot{\theta}, \dot{\psi})$

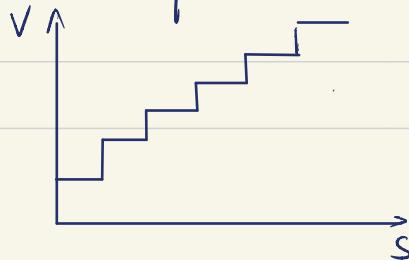
roll pitch yaw

2. Solution 1: Discretization



Problems:

① Naive representation for  $V^*$  (and  $\pi^*$ )

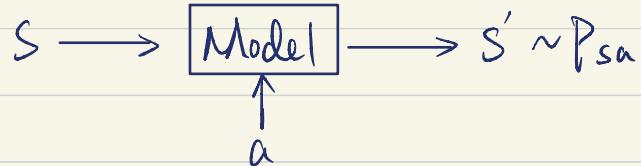


## ② Curse of Dimensionality

$S = \mathbb{R}^n$ , and discretize each dimension into  $k$  values, get  $k^n$  discrete states.

If dimension is low (2-6), can choose to use discretization.

## 3. Model (simulator) of MDP:



### ① Physics simulator

### ② Learn model from data: model-based RL

$$S_0^{(1)} \xrightarrow{a_0^{(1)}} S_1^{(1)} \xrightarrow{a_1^{(1)}} S_2^{(1)} \xrightarrow{a_2^{(1)}} \dots S_T^{(1)}$$

$$S_0^{(2)} \xrightarrow{a_0^{(2)}} \dots$$

$$S_0^{(m)} \xrightarrow{a_0^{(m)}} \dots$$

Apply supervised learning to estimate  $S_{t+1}$  as function of  $S_t, A_t$

e.g. linear regression:  $S_{t+1} = AS_t + BA_t$

$$\underset{A, B}{\text{Min}} \sum_{i=1}^m \sum_{t=0}^i \|S_{t+1}^{(i)} - (AS_t^{(i)} + BA_t^{(i)})\|^2$$

Model:

$$S_{t+1} = AS_t + BA_t \quad \leftarrow \text{deterministic}$$

or

works well in simulator, but usually works bad in reality.

$$\checkmark S_{t+1} = AS_t + BA_t + \varepsilon_t, \varepsilon_t \sim N(0, \sigma^2 I) \quad \leftarrow \text{stochastic}$$

#### 4. Fitted Value Iteration

Choose feature  $\phi(s)$  of state  $s$ .  $V(s) = \theta^\top \phi(s)$

Previously Value Iteration:

$$V(s) = R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V(s')$$

$$= R(s) + \gamma \max_a E_{s' \sim p_{sa}} [V(s')]$$

$$y^{(i)} \leftarrow \boxed{\max_a E_{s' \sim p_{sa}} [R(s) + \gamma V(s')]} \xrightarrow{g(a)}$$

Sample  $\{s^{(1)}, s^{(2)}, \dots, s^{(m)}\} \subseteq S$  randomly

Initialize  $\theta := 0$

Repeat {

For  $i = 1, \dots, m$  {

For each action  $a \in A$  {

Sample  $s'_1, s'_2, \dots, s'_k \sim P_{s^{(i)}} a$

Set  $g(a) = \frac{1}{k} \sum_{j=1}^k [R(s^{(i)}) + \gamma V(s'_j)]$

}

Set  $y^{(i)} = \max_a g(a)$

}

$\theta := \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^m (\theta^T \phi(s^{(i)}) - y^{(i)})^2$

Fitted VI gives approximation to  $V^*$ ,

Implicitly define  $\pi^*$ :  $\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{s' \sim P_{sa}} [V^*(s')]$ , use  $s'_1, s'_2, \dots, s'_k \sim P_{sa}$  to approximate expectation.

using the simulated model.

$\theta^T \phi(s'_j)$

last iteration of  $\theta$

\* Say model is  $S_{t+1} = f(S_t, A_t) + \xi_t$

For deployment (run-time): set  $\xi_t = 0$ , and  $k=1$

⇒ When in state  $s$ , pick action  $a$   $\underset{a}{\operatorname{argmax}} V(f(s,a))$   $s' \sim p_{sa}$   
simulation without noise