

## Reinforcement Learning 2

Review: MDP ( $S, A, \{P_{sa}\}, R, \mathbb{R}$ )

policy  $\pi: S \rightarrow A$

$V^\pi(s) \triangleq$  Expected total future payoff of policy  $\pi$  starting from  $s$

$$V^*(s) \triangleq \max_{\pi} V^\pi(s)$$

$\pi^*$  optimal policy

$\pi^*$  is greedy w.r.t.  $V^*(s)$

$$\pi^*(s) = \operatorname{argmax}_a (R(s) + \gamma \sum_{s'} P_{sa}(s') V^*(s'))$$

$$V^*(s) = V^{\pi^*}(s)$$

### Bellman Equations

linear system of equations  $V^\pi(s) = R(s) + \gamma \sum_{s'} P_{\pi(s)}(s') V^\pi(s')$

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

Value iteration, Policy iteration  
(VI) (PI)

VI, PI: known dynamics  $\{P_{sa}\}$

- not true in many real settings e.g. robotics

Assumption: Given  $s, a$  we can sample  $s' \sim P_{sa}$

model-based RL

vs.

model-free RL

explicitly learn transition dynamics from trajectories sampled

doesn't learn the transitions explicitly  
policy gradient (REINFORCE)  
 $\alpha$ -Learning

Tabular case

discrete state space

$$|S| < \infty, |A| < \infty$$

Continuous case

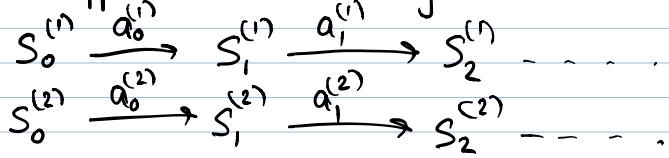
$$S = \mathbb{R}^d$$

$$A = \mathbb{R}^k$$

Model-based RL, tabular case

learn the dynamics

Suppose we have trajectories :



$$P_{sa}(s') = \frac{\text{\# times we arrived at } s' \text{ after taking action } a \text{ in state } s}{\text{\# times we took action } a \text{ in state } s}$$

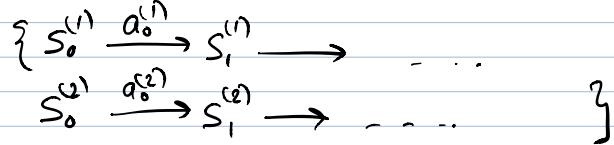
No Laplace smoothing

$$\text{replace } \frac{0}{0} \text{ by } \frac{1}{|S|}$$

Model-based RL algo:

Loop {  
  Initialize  $\pi$  randomly, dataset of trajectories  $D \leftarrow \emptyset$

(a) Execute  $\pi$  in the real environment to collect trajectories



add them to  $D$

(b) Estimate  $P_{sa}$  using data in  $D$

(Estimate the reward  $R$ )

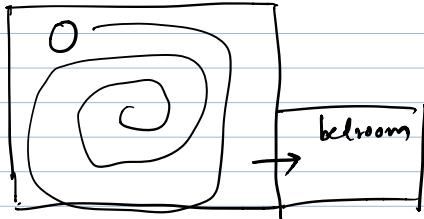
(c) invoke PI or VI to get  $V^*$  for estimated  $\{P_{sa}\}$

$$V^* = \arg \max_a \sum_{s'} P_{sa}(s') V^*(s')$$

$$\pi = \pi^*$$

Why loop?

- $P_{sa}$  in the beginning is not accurate
- initial policy may not see all the important states



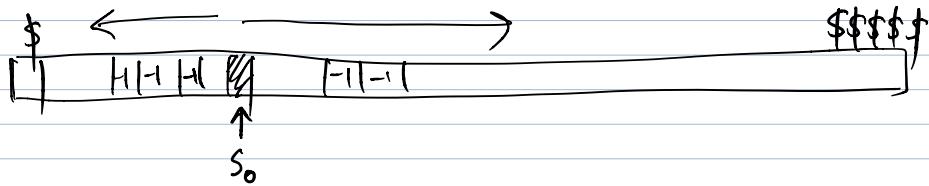
Exploitation vs.

Exploration trade off

find the best policy  
given the current  
information

collect more diverse data to get  
more information

exploration strategy: adding noise to actions generated  
by policy

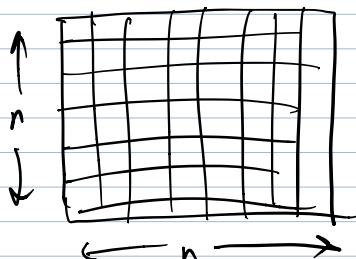


Continuous state space  $S = \mathbb{R}^d$

action space (discrete, continuous)

$d=2$  (or  $d=3$ )

discretization



$| \text{new state space} | = n^2$

$d > 4$  too many states!

Q1: how do we learn / represent the transition dynamics  $\{P_{sa}\}$

Q2: Value iteration (or PI)

A1 for Q1: physical simulator

A2: learning a model from data

represent the dynamics  $\{P_{sa}\}$

model  $s' \sim P_{sa}$  by

$$s' = f(s, a) + \xi$$

$f(\cdot)$  deterministic  
 $\xi$  noise  
(e.g. Gaussian)

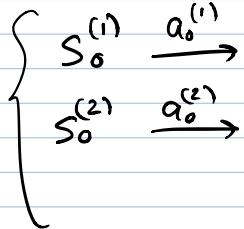
e.g.  $a \in \mathbb{R}^k$ ,  $s \in \mathbb{R}^d$

linear dynamics

$$s' = As + Ba + \xi$$

$$A \in \mathbb{R}^{d \times d}, B \in \mathbb{R}^{d \times k}$$

Linear Quadratic Regulator



⇒ supervised learning dataset

$$\{ (\underbrace{(s_0^{(1)}, a_0^{(1)})}_{x}, \underbrace{s_1^{(1)}}_{y}), \dots, (\underbrace{(s_{T-1}^{(1)}, a_{T-1}^{(1)})}_{x}, \underbrace{s_T^{(1)}}_{y}) \}$$

$$(\underbrace{(s, a)}_{x}, \underbrace{s'}_{y})$$

linear case: regression

$$A, B = \underset{\text{argmin}}{\sum_{i=1}^n \sum_{t=0}^{T-1} \|s_{t+1}^{(i)} - As_t^{(i)} - Ba_t^{(i)}\|_2^2}$$

non-linear

$f_\theta$  parametrized by  $\theta$

$$\theta = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n \sum_{t=0}^{T-1} \| s_{t+1}^{(i)} - f_\theta(s_t^{(i)}, a_t^{(i)}) \|_2^2$$

$$s' = f_\theta(s, a) + \xi$$

$$f_\theta(s, a) = A \phi_1(s) + B \phi_2(a) + \xi$$

$\uparrow$                      $\uparrow$   
feature map      feature map

$$f_\theta(s, a) = \text{neural network}(s, a)$$

Value iteration with continuous state space

Idea: Represent the value  $V(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$

$$V(s) = \theta^\top \phi(s) \quad \phi(s) : \text{feature map}$$
$$= \text{Neural net}(s; \theta)$$

Review: VI for discrete case

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

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

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

target value

Assume: A is discrete

Enforce BE on some randomly chosen  $S$   
 Loop {

(a) // estimate LHS of BE for randomly chosen  
 $s^{(1)}, \dots, s^{(n)}$

fix  $i$ ,  $s^{(i)}$  take  $s = s^{(i)}$

$\forall a \in A$  sample  $s_1', \dots, s_K' \sim p_{sa}$   $\leftarrow$  estimated dynamics

$$q(a) = R(s) + \frac{1}{K} \sum_{j=1}^K V(s_j')$$

estimator for  
 $E[V(s')]$   
 $s' \sim p_{sa}$

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

target value

(b) // enforce  $V(s^{(i)}) \approx y^{(i)}$

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

y

Fitted Value Iteration

$V^* \rightarrow \pi^*$  ?

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} p_{sa}(s') V^*(s')$$

$$= \underset{a}{\operatorname{argmax}} E_{s' \sim p_{sa}} [V^*(s')]$$

Choice 1: estimate the expectation by sampling

$$\approx \frac{1}{K} \sum_{i=1}^K V(s_i') \quad s_1', \dots, s_K' \sim p_{sa}$$

$$\begin{aligned} \text{Choice 2: } E_{s' \sim p_{sa}} [r^*(s')] &= E[V^*(f(s, a) + \xi)] \\ &\approx V^*(f(s, a)) \end{aligned}$$

Continuous A

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

$$y^{(i)} = \max_a (R(s^{(i)}) + \gamma \mathbb{E} V(f(s^{(i)}, a)))$$

by gradient ascent over actions