

## Reinforcement Learning

RL vs. Supervised learning

### ① Sequential decision making

return vs. risk tradeoff

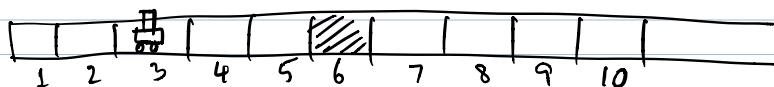
- long term future consequences
- greedy is not optimal

② no or little human supervision available

③ learn from designed reward fn

④ can collect more data iteratively

## Markov Decision Process



S: set of states (all possible configurations)

A: set of actions ( $A = \{L, R\}$ )

$P_{sa}$  dynamics / transitions

State transition prob. distrib'

At state  $s$ , when applying action  $a$

the dist. of next  $s'$

$$P_{sa}(s') = \Pr[s' | s, a]$$

$P_{sa}$  is a dist. ,  $P_{sa} \in \mathbb{R}^{|S|}$

$$\sum_{s' \in S} P_{sa}(s') = 1$$

Example: L action succeeds w. prob 0.9

$$P_{7,L}(6) = 0.9$$

$$P_{7,L}(7) = 0.1$$

$$P_{7,L}(s') = 0 \quad \text{if } s' \neq 6, 7$$

R action succeeds w. prob. 0.8

$$P_{7,R}(8) = 0.8$$

$$P_{7,R}(7) = 0.2$$

$$P_{7,R}(s') = 0 \quad \text{if } s' \neq 7, 8$$

Once you reach 6, robot stays there

$$P_{6,L}(6) = 1$$

$$P_{6,L}(s') = 0 \quad \text{if } s' \neq 6$$

$$P_{6,R}(6) = 1$$

$$P_{6,R}(s') = 0 \quad \text{if } s' \neq 6$$

Sequential decision making process

- $s_0$  initial state (given)
- algo chooses action  $a_0 \in A$
- environment step:  $s_1 \sim P_{s_0, a_0}$
- algo sees  $s_1$ , takes action  $a_1$
- environment:  $s_2 \sim P_{s_1, a_1}$

$$s_0 = 4$$

$$a_0 = R$$

$$s_1 \sim P_{4,R}$$

$$s_1 = 5$$

$$a_1 = R$$

Reward function

$$R: S \rightarrow \mathbb{R}$$

$$R(s)$$

$$R(6) = 1.0$$

$$R(s) = -0.1 \quad \text{if } s \neq 6$$

$$S_0 \xrightarrow{a_0} S_1 \xrightarrow{a_1} \dots$$

trajectory / episode

Total payoff:

$$R(S_0) + R(S_1) + \dots + R(S_t)$$

**Discounted total Payoff**

Discount factor:  $\gamma < 1$

$$R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots + \gamma^t R(S_t) + \dots$$

Suppose  $R(s) \in [-M, M]$

$\Rightarrow$  discounted total payoff  $\leq M + \gamma M + \gamma^2 M + \dots$

$$\frac{-M}{1-\gamma} \leq = \frac{M}{1-\gamma}$$

finite horizon: stop process at time  $T$

infinite horizon

Expected payoff

$$E[R(S_0) + \gamma R(S_1) + \dots + \gamma^t R(S_t) + \dots]$$

Reward can be  $R(s, a)$

$$MDP(S, A, \{P_{sa}\}_{\substack{s \in S \\ a \in A}}, \gamma, R)$$

goal: maximize expected discounted payoff

Markov property:

at any time  $t$ , given  $S_t$ , all future states  $S_{t+j}$

don't depend on past states

$$S_0, a_0, S_1, a_1, \dots, S_{t-1}, a_{t-1}$$

$\Rightarrow$  optimal action at time  $t$  does not depend on past states and actions

$\Rightarrow$  Best strategy can be expressed as a policy

$$\pi: S \rightarrow A \quad \text{policy}$$

$$a_t = \pi(s_t)$$

optimal policy

$$\pi(s) = R \quad \text{if } s \leq 6$$

$$\pi(s) = L \quad \text{if } s > 7$$

randomized policy  $\pi(a|s)$

$\exists$  deterministic optimum policy

Value function:

$$V^\pi: S \rightarrow \mathbb{R}$$

$V^\pi(s) \triangleq$  Expected payoff of executing policy  $\pi$   
Starting from state  $s$

$$s_0, a_0, \dots, a_t = \pi(s_t)$$

$$V^\pi(s) = \mathbb{E}[R(s_0) + \gamma R(s_1) + \dots + \gamma^t R(s_t) \mid s_0 = s]$$

How do we compute  $V^\pi(s)$



$$V^\pi(6) = 1 + \gamma + \gamma^2 + \gamma^3 + \dots$$

$$= \frac{1}{1-\gamma}$$

$$V^\pi(6) = 1 + \gamma (1 + \gamma + \gamma^2 + \dots)$$

$$= 1 + \gamma V^\pi(6)$$

$$V^\pi(6) = \frac{1}{1-\gamma}$$

$$V^\pi(s) = \mathbb{E}[R(s_0) + \gamma R(s_1) + \dots \mid s_0 = s]$$

$$= R(s) + \gamma \underbrace{E[R(s_1) + \gamma R(s_2) + \dots | s_0=s]}_{\text{total payoff starting from } s, \text{ and executing } \pi}$$

$$= R(s) + \gamma E[R(s_1) + \gamma R(s_2) + \dots]$$

$$s_i \sim P_{s, \pi(s)}$$

$$= R(s) + \gamma E[V^\pi(s_1)]$$

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{s, \pi(s)}(s') \cdot V^\pi(s')$$

### Bellman Equation

$V^\pi(s)$  as variables

$|S|$  variables,  $|S|$  equations

Solving system of linear eqn  $\Rightarrow V^\pi(s) \forall s$

given  $\pi, P_{s,a} \xrightarrow[\substack{\text{solve} \\ \text{linear system}}]{} V^\pi(s), \forall s$

Next:  $V^*(s) = \max_\pi V^\pi(s)$

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

Markov property  
 $\Rightarrow \pi^*$  that does not depend on  $s$   
 $\pi^*: S \rightarrow A$

### Bellman equation for $V^*$

$$V^*(s) = \max_\pi V^\pi(s)$$

$$= \max_\pi \left( R(s) + \gamma \sum_{s'} P_{s, \pi(s)}(s') \cdot V^\pi(s') \right)$$

$$= R(s) + \max_\pi \left( \gamma \sum_{s'} P_{s, \pi(s)}(s') \cdot V^\pi(s') \right)$$

$$= R(s) + \max_a \left( \gamma \sum_{s'} P_{sa}(s') \max_\pi V^\pi(s') \right)$$

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

the best expected payoff  
after taking action  $a$

$$V^* \in \mathbb{R}^{|S|}$$

$$\begin{bmatrix} V^*(1) \\ \vdots \\ V^*(|S|) \end{bmatrix}$$

$$V^* = B(V^*)$$

$$B \text{ operator} \quad \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$$

### Value iteration

- Initialize  $V \in \mathbb{R}^{|S|}$ ,  $V=0$   $V(s)=0 \forall s$

Loop

$$V = B(V)$$

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

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s' \in S} P_{sa}(s') \cdot V^*(s')$$

### Policy iteration

- initialize policy  $\pi$

- loop {

- let  $V = V^\pi$  (obtained by solving system of eqn)

$$\pi(s) = \operatorname{argmax}_a \sum_{s'} P_{sa}(s') V(s')$$

}

$V^*$ ,  $\pi^*$  is stationary point of algorithm

$$V^* = V^{\pi^*}$$

Common:  $V^*$ ,  $\pi^*$  are stationary points of algos.

### Policy iteration:

Con: each iteration requires solving linear system  
expensive for large state space

Pro: converges in finite time exactly

At most  $|A|^{st}$  policies

Value iteration always has some approximation

### Practice: discrete state space

- value iteration often slightly better
- continuous controls policy iteration  
(or some combina<sup>n</sup>)

because max operation in Value iteration is hard to do exactly