

# Reinforcement Learning

RL vs Supervised learning

① Sequential decision making

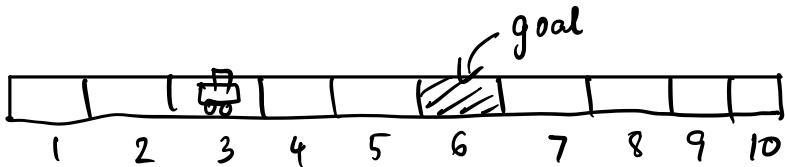
- long term future consequences
- greedy is not optimal
- return vs. risk tradeoff

② No or little human supervision

③ Learn from designed reward function

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 , when applying action a  
the dist' of next  $s'$

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

$P_{sa}$  is a dist'n'

$$P_{sa} \in \mathbb{R}^{|S|} \quad \sum_{s' \in S} P_{sa}(s') = 1$$

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

$$P_{7,R}(8) = 0.8 \quad P_{7,R}(7) = 0.2$$

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

Once you reach 6, robot stays there:

$$P_{6,L}(6) = 1, \quad P_{6,L}(s') = 0 \quad \forall s' \neq 6$$

$$P_{6,R}(6) = 1, \quad \dots$$

## Sequential decision making process

- $s_0$  initial state (given)
- algo chooses action  $a_0 \in A$
- environment step:  $s_1 \sim P_{s_0, a_0}$  (Partially observed MDP)
- 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

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

$$R(s)$$

$$R(6) = 1.0$$

$$R(s) = -0.1 \quad \forall s \neq 6$$



Total payoff:

$$R(s_0) + R(s_1) + \dots + R(s_t) + \dots$$

Discounted total payoff (Infinite horizon)

Discount factor  $\gamma < 1$  (interest rate)

$$R(s_0) + \gamma R(s_1) + \dots + \gamma^t R(s_t) + \dots$$

Discounted reward is bounded

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

$\Rightarrow$  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  $T$

Expected payoff:

$$\mathbb{E}[R(s_0) + \gamma R(s_1) + \dots + \gamma^t R(s_t) + \dots]$$

Reward can be  $R(s, a) \leftarrow$  function of state & action

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

goal: maximize expected discounted payoff.

Markov property:

at any time  $t$ , given  $s_t$ , all future states

$s_{t+j}$  doesn't depend on past states/actions

$s_0, a_0, \dots, s_{t-1}, a_{t-1}$

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

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

$\pi: S \rightarrow A$  policy

$a_t = \pi(s_t)$

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

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

randomized policy  $\pi(a|s)$

$\exists$  deterministic optimal policy

Value function :

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

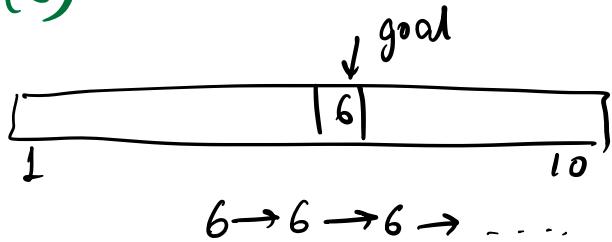
$V^\pi(s) \triangleq$  Expected total payoff of executing policy  $\Pi$  starting from state  $s$

$s_0, a_0, \dots$

$$a_t = \pi(s_t)$$

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

How to compute  $V^\pi(s)$



$$\begin{aligned} V^\pi(6) &= 1 + \gamma + \gamma^2 + \dots + \gamma^t + \dots \\ &= \frac{1}{1-\gamma} \end{aligned}$$

$$V^\pi(6) = 1 + \gamma \underbrace{(1 + \gamma + \gamma^2 + \dots)}_{V^\pi(6)}$$

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

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

$$\begin{aligned} V^\pi(s) &= \mathbb{E}[R(s_0) + \gamma R(s_1) + \dots \mid s_0 = s] \\ &= R(s) + \gamma \underbrace{\mathbb{E}[R(s_1) + \gamma R(s_2) + \dots \mid s_0 = s]}_{\text{total payoff starting from } s \text{ and executing } \Pi} \\ &= R(s) + \gamma \mathbb{E}_{s_1 \sim P_{s, \pi(s)}} [R(s_1) + \gamma R(s_2) + \dots] \\ &= R(s) + \gamma \mathbb{E}_{s_1 \sim P_{s, \pi(s)}} [V^\pi(s_1)] \\ &= R(s) + \gamma \sum_{s' \in S} P_{s, \pi(s)}(s') \cdot V^\pi(s') \end{aligned}$$

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

Bellman  
Equation

$V_{\pi}(s)$  as variables

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

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

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

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

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

Markov Property  $\Rightarrow \exists \pi^*$  that simultaneously maximizes  $V^{\pi}(s) \forall s \in S$

$$\pi^*: s \rightarrow A$$

Instead of computing  $\pi^*$  directly

$$\text{let's calculate } V^*(s) \triangleq V^{\pi^*}(s)$$

We will get  $\pi^*$  from  $V^*(s)$

Bellman Equation for  $V^*$

$$\begin{aligned} V^*(s) &= \max_{\pi} V^{\pi}(s) \\ &= \max_{\pi} \left( R(s) + \gamma \sum_{s'} P_{s, \pi(s)}(s') \cdot V^{\pi}(s') \right) \\ &= R(s) + \max_{\pi} \left( \gamma \sum_{s'} P_{s, \pi(s)}(s') \cdot V^{\pi}(s') \right) \\ &= R(s) + \max_a \left( \gamma \sum_{s'} P_{s, a}(s') \cdot \max_{\pi} V^{\pi}(s') \right) \\ &= R(s) + \gamma \max_a \sum_{s'} P_{s, a}(s') \cdot V^*(s') \end{aligned}$$

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

Bellman Operator

best expected payoff  
after taking action a

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

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

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

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

## Value Iteration

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

$$V := B(V)$$

$\uparrow$  new       $\uparrow$  old

$$\forall s \in S, V(s) = R(s) + \gamma \max_a \sum_{s'} P_{s,a}(s') \cdot V(s')$$

$V \rightarrow V^*$  because  $B$  is a contraction

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

## Policy Iteration

Initialize a policy  $\pi$

- Loop {

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

$$- \pi(s) = \operatorname{argmax}_a \sum_{s'} P_{s,a}(s') \cdot V(s')$$

}

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

## Policy Iteration

Con: each iteration needs solving linear system

Pro: converge in finite time complexity  
 $A$  is finite  
 $S$  is finite  
At most  $|A|^{|\mathcal{S}|}$  policies

Practice: discrete state space

- Value iteration often slightly better
- continuous controls : policy iteration