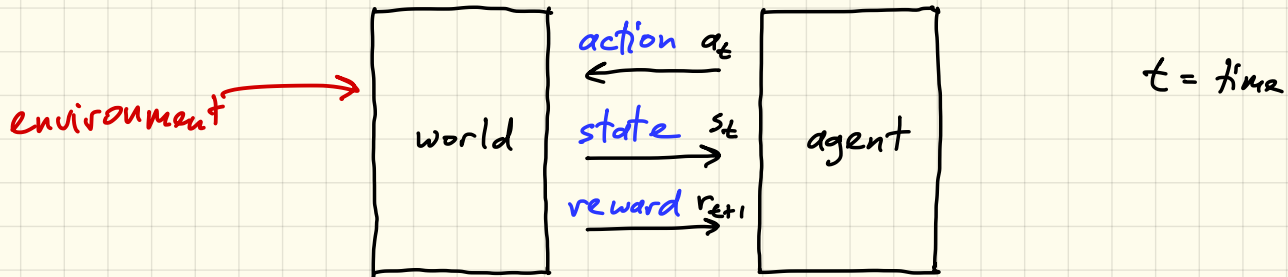


MDP: Markov Decision Process (S&B ch 3)

- stochastic, discrete state, discrete action, state feedback



3 system spaces: state $s_t \in S$, action $a_t \in A$, reward $r_t \in \mathbb{R}$

3 system functions: model $p(s_{t+1} | s_t, a_t) = \text{probability of } s_{t+1} \text{ given } s_t, a_t$

policy $\pi(a_t | s_t) = \text{probability of } a_t \text{ given } s_t$

reward $r_{t+1} = R(s_t, a_t) \leftarrow \text{not stochastic (if so, } E[R])$

Types of ML

supervised learning:

given many pairs (x_i, y_i)

instruction

learn approximate map $y = f_\theta(x)$

$$y = \theta_0 + \theta_1 x$$

$$\text{minimize}_{\theta} \sum_i \|y_i - f_\theta(x_i)\|^2$$

unsupervised learning:

given many x_i

- clustering, k-means
- autoencoders
- GANs (generative adversarial networks)

self-supervised learning:

given sequence $(x_0, x_1, x_2, \dots, x_n)$

learn model $x_{k+1} = f_\theta(x_k, x_{k-1}, x_{k-2}, \dots, x_{k-m})$

- OpenAI GPT-3

reinforcement learning : $a_t = \pi_{\theta}(s_t)$ $(s_{t+1}, r_t) = p(s_t, a_t)$

policy (arrow from π_{θ} to a_t)

recursive (arrow from p to s_{t+1})

evaluation (arrow from p to r_t)

$$\underset{\theta}{\text{maximize}} \sum_t r_t$$

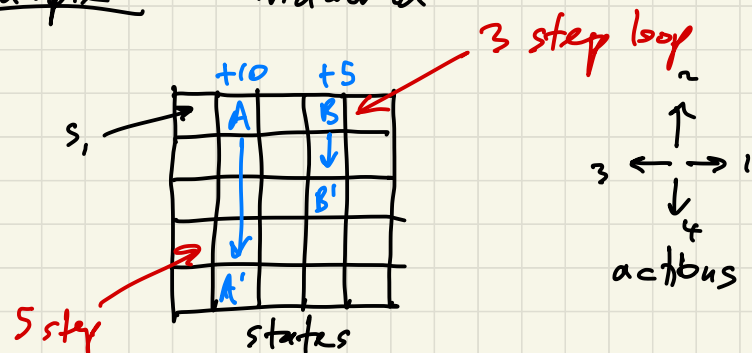
model-based RL : algorithms can directly call $p()$

model-free RL : only have access to $p()$ via trajectories

↳ pros : - can do it on reality
- only need simulator

can we
call $\nabla p()$ (arrow from $\nabla p()$ to $p()$ in the model-based RL text)

Example : Gridworld



5 step loop

$$s \in \{1, \dots, 25\}$$

(0 to 24 in Python)

$$a \in \{1, \dots, 4\}$$

reward: $r_t = -1$ except
at A, B

If the system is time-dependent, we can add time as a state.

Value functions

$$V_{\pi}(s) = E \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \pi, S_0 = s \right]$$

← state value function

← infinite horizon

← policy π

← discount factor $\gamma < 1$

← expected total future discounted reward starting from s under policy π

effective time horizon $\sim \frac{1}{\gamma}$

smaller $\gamma \Rightarrow$ prioritize short-term rewards

typical $\gamma = 0.95, 0.99$

(or use finite horizon $V_{\pi}(s) = E \left[\sum_{t=0}^T \gamma^t r_{t+1} \right]$)

← finite horizon

$$Q_{\pi}(s, a) = E \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \pi, S_0 = s, A_0 = a \right]$$

← state-action value function

Relationships :

def: $V_{\pi}(s) = Q_{\pi}(s, \pi(s))$

$$V_{\pi}(s) = E_{a \sim \pi(\cdot|s)} [Q_{\pi}(s, a)] = \sum_a \pi(a|s) Q_{\pi}(s, a)$$

$$Q_{\pi}(s, a) = R(s, a) + \gamma E_{s' \sim p(\cdot|s, a)} [V_{\pi}(s')]$$

$$= R(s, a) + \gamma \sum_{s'} p(s'|s, a) V_{\pi}(s')$$

def: $Q_{\pi}(s, a) = R(s, a) + \gamma V_{\pi}(\underbrace{p(s, a)}_{s'})$

recursive relation for V :

$$V_{\pi}(s) = E_{a \sim \pi(\cdot|s)} [R(s, a) + \gamma E_{s' \sim p(\cdot|s, a)} [V_{\pi}(s')]]$$

Optimality: $V^*(s) = \max_{\pi} V_{\pi}(s)$

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} Q^*(s, a') \\ 0 & \text{otherwise} \end{cases}$$