

Action affects world state

- choose classes

- action: take a class
- rewards: enjoyment
- state: satisfy more prefs

e.g. a class may be good to take later,
bad to take now

- robotics

- action: motor
- reward: complete a task
- state: position of robot

state transitions are noisy/random

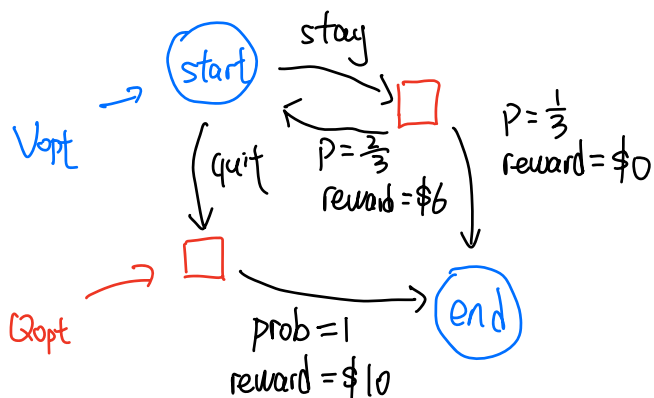
- video games

Markov Decision Process (MDP)

Formal description of a world with state actions, rewards, etc.

At each time:

- player can stay or quit
- if quit: get \$10 and game ends
- if stay:
 - prob $1/3$, get \$0 and end
 - prob $2/3$, get \$6 and continuous



Formal MDP:

- set of states (e.g. possible positions of robots)
- starting state: S_{start}
- Actions(s): possible actions at states s
- $T(s, a, s')$: prob of going from state s to s' after taking action a .
(e.g. $T(start, stay, end) = \frac{1}{3}$)
- $R(s, a, s')$: reward of going from s to s' taking action a .
- $IsEnd(s)$: is this an end state

Give MDP, what is optimal agent behavior?

policy: A strategy that agent can use

$$\underbrace{\pi(s)}_{\text{current state}} \rightarrow \underbrace{\text{action} \in \text{Action}(s)}_{\text{chosen action}} \quad \text{discounted}$$

The value $V_{\pi}(s)$ for policy π is expected \checkmark sum of rewards starting from state s , applying policy π .

Discounting: future rewards are less valuable than current rewards

- at any timestep you could die

we introduce a discount factor $\gamma \in [0, 1]$
prob of survival at each timestep

If we get a sequence of rewards r_1, r_2, \dots

Discounted sum of rewards = $r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$

The optimal value $V_{opt}(s)$ is maximum possible value at state s for any policy.

V_{opt} is characterized by recursive formulas:

$$V_{opt}(s) = \begin{cases} 0 & \text{if } isEnd(s) \\ \max_{\alpha \in Actions(s)} Q_{opt}(s, \alpha) & \text{else} \end{cases}$$

expected optimal value after taking action α in state s

$$Q_{opt}(s, \alpha) = \sum_{s'} T(s, \alpha, s') [R(s, \alpha, s') + \gamma V_{opt}(s')]$$

prob of transition to s' reward now discounted future rewards at s'

optimal policy: $\pi^*(s) = \operatorname{argmax}_{\alpha \in Action(s)} Q_{opt}(s, \alpha)$

if we can estimate $Q_{opt}(s, \alpha)$ for all s, α , we can find π^* .

RL:

- believe the world is a MDP
- don't know $T(s, a, s')$ and $R(s, a, s')$
- has to try many actions in many states

For episode = 1, 2, 3 ...

$S_1 \leftarrow S_{\text{start}}$

for $t=1, 2, \dots$

- agent chooses action $a_t = \pi_{\text{act}}(s_t)$

policy we act with during learning

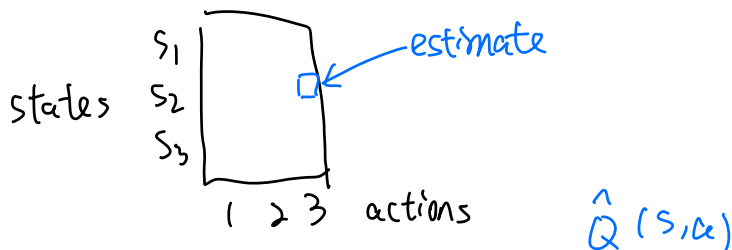
- agent receives :

- reward r_t

- new state S_{t+1}

- Update agent's parameters

Q-learning : directly learn $Q_{\text{opt}}(s, a)$



$$Q_{\text{opt}}(s, a) = \sum_{s'} T(s, a, s') \cdot [R(s, a, s') + \gamma V(s')]$$

$$V_{\text{opt}}(s') = \begin{cases} 0 & \text{if } \text{IsEnd}(s') \\ \max_{a \in \text{Action}(s')} Q_{\text{opt}}(s', a) \end{cases}$$

We have data

$$\underbrace{s_1, a_1, r_1, s_2, a_2, r_2, \dots}_{1 \text{ example}} \quad \underbrace{\hspace{10em}}_{2 \text{ example}}$$

Every time we see (s, a, r, s') :

"nudge" $\hat{Q}(s, a)$ based on this observation

$$\hat{Q}(s, a) \leftarrow (1 - \eta) \hat{Q}(s, a) + \eta (r + \gamma \hat{V}(s'))$$

↑
↑
↑
↑

estimate of
learning rate
observed reward
observed next state

$Q_{\text{opt}}(s, a)$
e.g. 0.1
one sample for $Q_{\text{opt}}(s, a)$

$$\text{where } \hat{V}(s') = \begin{cases} 0 & \text{if } \text{IsEnd}(s') \\ \max_{a \in \text{Actions}(s')} \hat{Q}(s', a) \end{cases}$$

What π_{act} to act with?

$$\pi(s) = \arg \max_{a \in \text{Actions}(s)} \hat{Q}(s, a)$$

pure exploitation strategy

solution: ϵ -greedy

policy during learning $\pi_{\text{Act}} = \begin{cases} \text{with prob } 1-\epsilon, \text{ do } \operatorname{argmax}_a \hat{Q}(s,a) \\ \text{with prob } \epsilon, \text{ do random action } \in \text{Actions}(s) \end{cases}$

- Training: Q-learning with $\epsilon = 0.1$

balance exploration v.s. exploitation

- Testing: act with $\epsilon = 0$