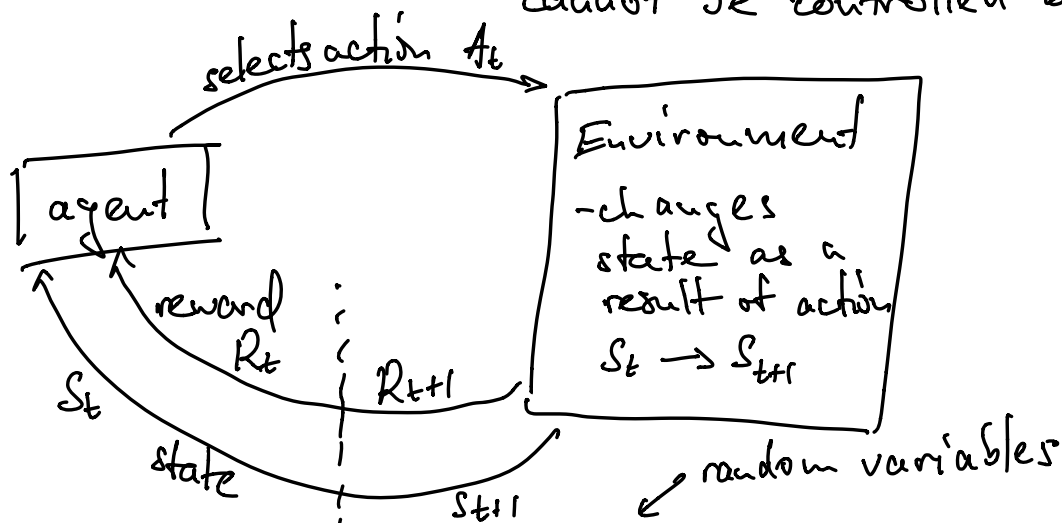


Markov Decision Processes (MDPs)

- Agent - Environment Interface

agent: learner / decision maker

environment (env.): everything outside agent;
cannot be controlled by agent



defs.: 1) state space: $S \ni S_t$

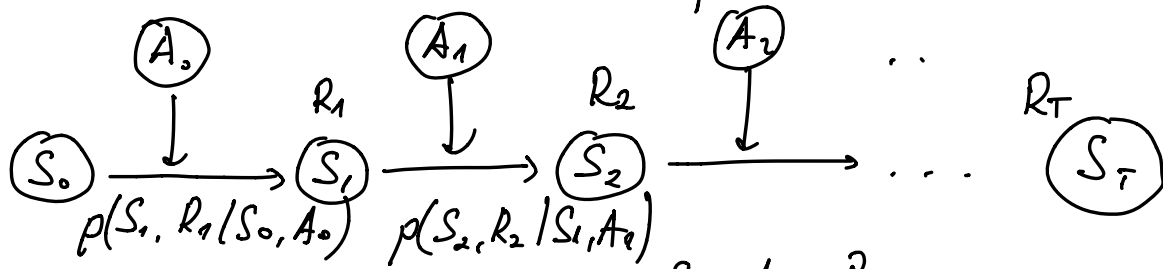
2) action space: $A \ni A_t$ / allowed actions may depend on state

3) reward space: $R \ni R_t$

- trajectory: $S_0 A_0 R_1 \underbrace{S_1 A_1}_{t=1} R_2 S_2 \dots$
step: $t=0 \quad \quad \quad t=1 \quad \quad \dots$

→ agent finds env. in the state S_0 ,
chooses action A_0 ,
which causes the env. to transition
to state S_1 ,
and the agent receives a reward R_1
as a consequence

- finite Markov decision process (MDP)



- a Markov process on $S \times A \times R$ space
 - provide a mathematical abstraction for the problem: goal-directed learning from interactions
 - the transition prob. p depends only on the preceding state and action, not on the history
- $$p(s', r | s, a) := \Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$
- for a $s, s' \in S, a \in \mathcal{A}(s), r \in R$
- defines the dynamics of the MDP

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \quad \forall s \in S, a \in \mathcal{A}(s)$$

- state transition probability

$$p(s' | s, a) := \Pr(S_t = s' | S_{t-1} = s, A_{t-1} = a)$$

$$= \sum_{r \in R} p(s', r | s, a)$$

(marginalize/integrate out rewards)

→ if rewards are stochastic:

$$r(s,a) := \mathbb{E}_p [R_t / S_{t-1} = s, A_{t-1} = a]$$

$$= \sum_{r \in R} \sum_{s' \in S} r p(s', r | s, a)$$

expected reward for state-action pair (s,a)

$$r(s,a,s') := \mathbb{E}_p [R_t / S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$$= \frac{\sum_{r \in R} r p(s', r | s, a)}{\sum_{r \in R} p(s', r | s, a)} = \sum_{r \in R} r \frac{p(s', r | s, a)}{p(s' | s, a)}$$

— Example: (recycling robot)

goal: teach a robot to collect empty soda cans

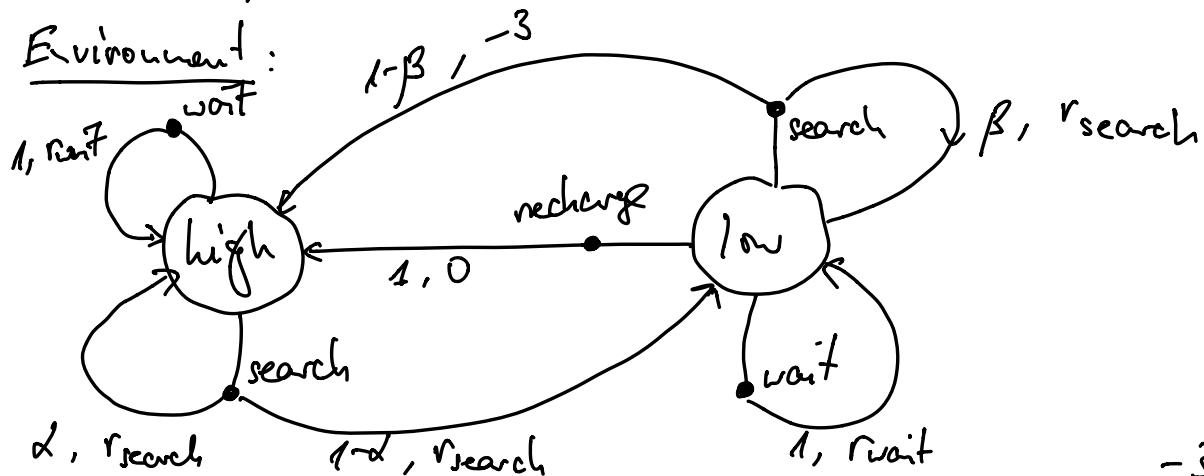
• state space: $S = \{\text{low}, \text{high}\}$ models battery level

• action space: $A(\text{high}) = \{\text{search}, \text{wait}\}$

$A(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}$

• reward space: $R = \{r_{\text{search}}, r_{\text{wait}}, 0, -3\}$

Environment:



s	a	s'	$p(s'/s,a)$	$r(s,a,s')$
high	search	high	α	r_{search}
high	search	low	$1-\alpha$	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	-
low	search	high	$1-\beta$	-3
low	search	low	β	r_{search}
low	wait	low	1	r_{wait}
low	wait	high	0	-
low	recharge	low	0	-
low	recharge	high	1	0

-rewards: $R_t \in \mathbb{R} \subset \mathbb{R}$

- communicate to agent what it has to achieve, not how to achieve it

- instantaneous

- goal for agent: maximize total expected return

→ sequence of rewards within trajectory:

$$R_1, R_2, R_3, \dots, R_T$$

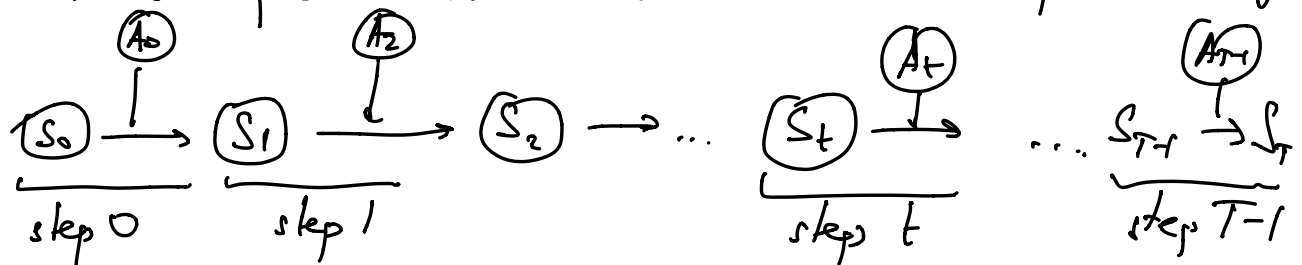
return: $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$
reward to-go

— two kind of tasks:

→ episodic / finite-horizon tasks:

• agent learns in multiple repeated trials called episodes (e.g. plays of a game)

→ each episode is broken down in steps labelled by t



→ subsequent episodes begin independently of how current episode ended

→ return: $G_t := \sum_{k=t+1}^T R_k$

→ non-episodic / infinite-horizon tasks

— no episodes, continue forever

problem: $G_t = \infty$

→ way out: introduce discounted rewards
agent cares more about near-term rewards!

$$G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

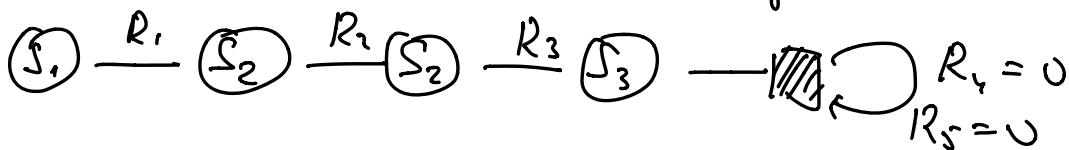
$$= \sum_{k=0}^{\infty} \gamma^k R_{t+1+k} \quad , \quad \gamma \in [0, 1] \text{ discount factor}$$

$\Rightarrow G_t < \infty$ if $\{R_k\}_k$ is bounded

eg. $R_k = 1 \ \forall k \Rightarrow G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$
 geometric series

— unified notation:

$$G_t := \sum_{k=0}^T \gamma^{k-t-1} R_k : \quad \begin{cases} T = \infty & : \text{non-episodic} \\ T < \infty & : \text{episodic} \\ \gamma = 1 \end{cases}$$



Examples :

1) finite-horizon tasks

→ recycling robot

→ solve a maze

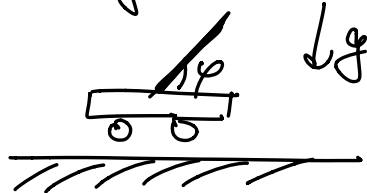
→ board games (chess, backgammon)

2) infinite-horizon tasks:

→ cart pole problem:

(pole-balancing)

→ teach a robot to walk (forever)



Policies & Value Functions

How does the RL agent choose actions?

idea: agent observes state S_t of the env.

Can we correlate states & expected returns?

value function: estimate "how good" it is for the env. to be in a given state

note: future rewards depend on what actions the agent will take

→ value functions are defined wrt particular ways of acting, called policies

policy: $\pi: \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$
 $(a, s) \mapsto \pi(a|s)$

probability to take action a in the state s

$$\sum_{a \in \mathcal{A}(s)} \pi(a|s) = 1 \quad (\text{normalized})$$

1) deterministic policy: ("delta"-function policy)

$$\pi(a|s) = \delta_{a,a_0} = \begin{cases} 1, & \text{if } a = a_0 \\ 0, & \text{otherwise} \end{cases}$$

2) stochastic policy has weight on multiple actions.

eg. $\begin{matrix} 30\% & 15\% & 15\% & 10\% & 30\% \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ a_0 & a_1 & a_2 & a_3 & a_4 \end{matrix}$

- expectation value of R_{t+1} under π

$$\mathbb{E}_{\pi}[R_{t+1} | S_t = s] = \sum_{a \in A(s)} \pi(a|s) \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) r$$

- value function (of a policy π)

$$V_{\pi}: S \rightarrow \mathbb{R} \\ s \mapsto V_{\pi}(s)$$

$$V_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right]$$

- expected return starting from a state s under the policy π

$\mathbb{E}_{\pi}[\cdot]$: expectation of a random variable given that agent follows policy π at every step t in the episode!

- recursive relation for $v_\pi(s)$

recall: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$
 $= R_{t+1} + \gamma \underbrace{(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots)}_{=: G_{t+1}}$

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi [G_t | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_\pi [R_{t+1} | S_t = s] + \gamma \mathbb{E}_\pi [G_{t+1} | S_t = s] \\ &= \sum_{a \in A(s)} \pi(a|s) \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) r \\ &\quad + \gamma \sum_{a \in A(s)} \pi(a|s) \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) \underbrace{\mathbb{E}_\pi [G_{t+1} | S_{t+1} = s']}_{= v_\pi(s')} \\ &= \sum_{a \in A(s)} \pi(a|s) \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) (r + \gamma v_\pi(s')) \end{aligned}$$

relates $v_\pi(s)$ & $v_\pi(s')$

→ Bellman equation for v_π :

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) [r + \gamma v_\pi(s')]$$

• averages all possibilities for trajectories
weighting each by its probability for occurring

- the value of the starting state $v(s)$ is equal to the discounted value of the expected next state, plus the reward expected along the trajectory

- v_π is the unique solution of the Bellman eq. (finite MDPs)

- Action-value functions : q -functions (of a policy π)

$$q_\pi: S \times A \rightarrow \mathbb{R}$$

$$(s, a) \mapsto q_\pi(s, a)$$

$$q_\pi(s, a) := \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a]$$

→ the expected return starting in a state s and taking the action a , and following the policy π afterwards

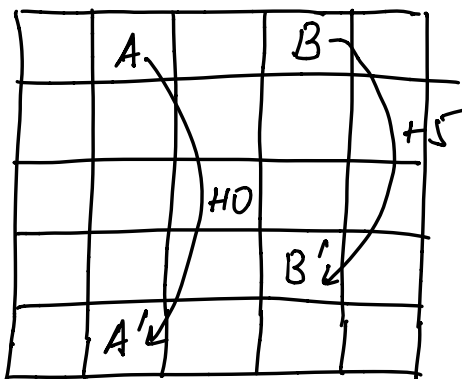
show (HW):

$$q_\pi(s, a) = \sum_{r \in \mathbb{R}} \sum_{s' \in S} p(s', r \mid s, a) [r + \gamma v_\pi(s')]$$

- relation between $v_\pi(s)$ & $q_\pi(s, a)$?

$$v_\pi(s) = \sum_{a \in A(s)} \pi(a \mid s) q_\pi(s, a)$$

Example: Gridworld: $r(s, a, s')$



rewards: $r(A, :, A') = +10 \quad \forall a \in A$
 $r(B, :, B') = +5 \quad \forall a \in A$
 $r(\text{boundary}, a_{\text{outside}}, s') = -1$
 $r(s, a, s') = 0$, otherwise

states: squares on grid

actions: \leftrightarrow ; cannot move across boundaries

policy: $\pi(a/s) = \frac{1}{|A|} = \frac{1}{4}$ equiprobable random policy

→ value function $v_{\pi}(s)$ for $\gamma = 0.9$

3.3	8.8	4.7	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Optimal policies & value functions

- ordering relation for policies:

$$\pi \geq \pi', \text{ iff } V_\pi(s) \geq V_{\pi'}(s) \quad \forall s \in S$$

there exists (at least one) policy that is better than or equal to all other policies

\Rightarrow optimal policy: π_*

- optimal value function: V_*

\rightarrow value function associated with π_*

$$V_*(s) := \max_{\pi} V_\pi(s) \quad \forall s \in S$$

- optimal q-function: q_*

$$q_*(s, a) := \max_{\pi} q_\pi(s, a) \quad \begin{array}{l} \forall s \in S \\ \forall a \in A \end{array}$$

$$\Rightarrow q_*(s, a) = \mathbb{E} [R_{t+1} + V_*(S_{t+1}) \mid S_t = s, A_t = a]$$

- relation:
$$V_*(s) = \max_{a \in A(s)} q_*(s, a)$$

\rightarrow note: optimal policy π_* may not be unique but V_* & q_* are unique

- Bellman optimality equation for $v_*(s)$:

$$v_*(s) = \max_{a \in A(s)} q_*(s, a) = \max_{a \in A(s)} q_{\pi_*}(s, a)$$

$$= \max_{a \in A(s)} \mathbb{E}_{\pi_*} [G_t \mid S_t = s, A_t = a]$$

$$= \max_{a \in A(s)} \mathbb{E}_{\pi_*} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$= \max_{a \in A(s)} \mathbb{E} [R_{t+1} + \gamma \underbrace{v_{\pi_*}(S_{t+1})}_{= v_*(S_{t+1})} \mid S_t = s, A_t = a]$$

$$= \max_{a \in A(s)} \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')]$$

→ if v_* is known, then an optimal policy can be defined by taking all possible actions from s , and comparing the $v_*(s')$ values

⇒ greedy policy (deterministic)

→ v_* turns the total expected return to a local quantity, available at each step.

⇒ a one-step-ahead search with v_* yields long-term optimal actions

- Bellman optimality equation for $q_*(s,a)$

$$\begin{aligned} q_*(s,a) &= E[R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= E[R_{t+1} + \gamma \max_{a' \in A(S_t)} q_*(S_{t+1}, a') \mid S_t = s, A_t = a] \\ &= \sum_{r, s'} p(s', r \mid s, a) [r + \gamma \max_{a' \in A(s')} q_*(s', a')] \end{aligned}$$

→ resolution over the actions a

$$\pi_*(a \mid s) = \underset{a \in A(s)}{\operatorname{argmax}} q_*(s, a) \quad \text{optimal greedy policy}$$

→ no need for a one-step-ahead search

- Remarks:

→ exact solution to Bellman optimality eq.'s nearly feasible, since

1) do not know transition prob.

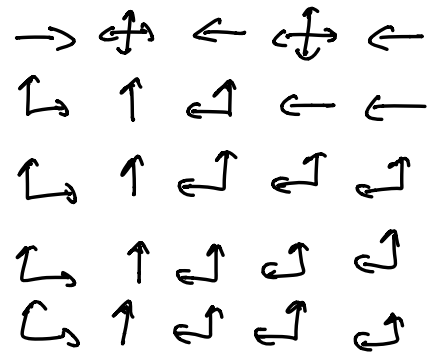
$p(s', r \mid s, a)$ (→ use MC sampling to find high-probability trajectories)

2) state space S is often exponentially large → cannot find $V_*(s)$, $q_*(s,a)$ & π_*

- Example : optimal policy for Gridworld:
($\gamma = 0.9$)

22.8	22.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

$V_*(s)$



π_*