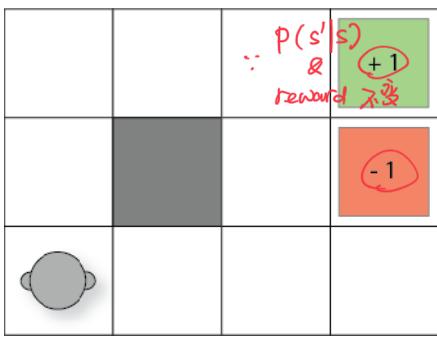


Probabilistic planning - Markov Decision Processes (MDPs)



- An agent has a goal to navigate cells
- The grey square is a wall (like the edges of grid)
 - The two coloured cells giving rewards: 1 (**goal**) and -1 (**bad goal**)

Actions have **non-deterministic** outcomes (effects)!

- If the agent tries to move north, 80% of the time, this works as planned (provided the wall is not in the way)
- 10% of the time, trying to move north takes the agent east (provided the wall is not in the way)
- 10% of the time, trying to move north takes the agent west (provided the wall is not in the way);
- If wall is in the way of the cell that would have been taken, the agent stays put
- Similar for all other directions

MDPs:

- Set of states S
- Initial state I
- Probabilistic state transitions: $\sum_{s'} P(a|s, s') = 1 \in [0, 1]$
- Reward function $r(s, a, s')$ in Real
- Discount factor γ (gamma)

Classical Planning:

- Set of states S
 - Initial state I
 - Transition function A deterministic
 - Goals G
 - Costs
- $$\begin{cases} S \xrightarrow{a} S' \\ S \xrightarrow{a=1} S' \\ S \xrightarrow{a=0} S' \end{cases}$$
- $\} MDP = \text{classical planning}$

γ discount factor is $(0, 1)$ not inclusive
Discounted rewards

Reward at T time

$$\begin{aligned} G_t &= r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4 \dots \\ &= r_1 + \gamma(r_2 + \gamma(r_3 + \gamma(r_4 \dots))) \\ &= r_t + \gamma G_{t+1} \end{aligned}$$

Modelling MDPs --- Probabilistic PDDL

```
(define (domain bomb-and-toilet)
  (:requirements :conditional-effects :probabilistic-effects)
  (:predicates (bomb-in-package ?pkg) (toilet-clogged)
    (bomb-defused))

  (:action dunk-package
    :parameters (?pkg)

    :effect (and (when (bomb-in-package ?pkg)
      (bomb-defused))
      (probabilistic 0.05 (toilet-clogged))))
```

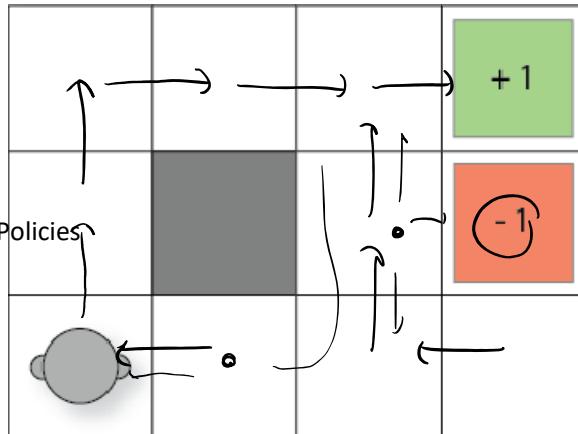
planning: action sequence

Mapping: states to actions

Solution for MDP is a *policy*:
function

$\text{at}(0,0) \Rightarrow \text{move_up}$
 $\text{at}(0,1) \Rightarrow \text{move_up}$
 $\text{at}(0,2) \Rightarrow \text{move_right}$
 $\text{at}(1,0) \Rightarrow \text{move_left}$
 $\text{at}(1,2) \Rightarrow \text{move_right}$
 $\text{at}(2,0) \Rightarrow \text{move_up}$
 $\text{at}(2,1) \Rightarrow \text{move_up}$
 $\text{at}(2,2) \Rightarrow \text{move_right}$
 $\text{at}(3,0) \Rightarrow \text{move_left}$

(solutions to MDPs -- Policies)



policies {

$\pi(s) \rightarrow A$	deterministic	same action in the same state → Subject focus
$\pi(s,a) \in [0,1]$	stochastic	diff action in the same state

Solving MDPs

Expected return exercise:

You can steal:

- A) An iPhone, which you think you have a 20% chance of selling for \$500, or an 80% chance of selling for \$250.
- B) A Samsung, which you think you have a 50% chance of selling for \$500, or a 50% chance of selling for \$200.

A: $0.2 \cdot 500 + 0.8 \cdot 250 = 300$

B: $0.5 \cdot 500 + 0.5 \cdot 200 = 350$

State (s)
of
value ⇒ what action is the best

Bellman equation: (Value func)

$$V(s) = \max_{a \in A} \sum_{s' \in S} P_a(s'|s) \left[r(s, a, s') + \gamma V(s') \right]$$

all the possible s' we can transition to

$\Delta(s, a)$
↑
quality

expected reward action a

immediate reward
future reward
Discount