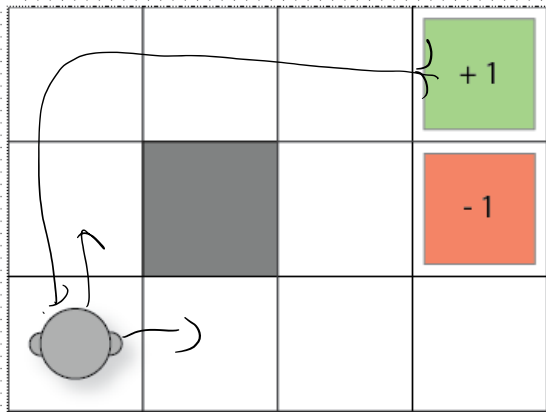


# Markov Decision Processes

Tuesday, 4 September 2018 10:29 AM



The grey square is a wall.

The two labelled cells give a reward: 1 and -1 respectively.

But! Things can go wrong:

- 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 west (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)
- If the wall is in the way of the cell that would have been taken, the agent stays put.
- Similar for all other directions

Classical Planning:

- Set of states  $S$
- Initial state  $I$
- Transition function  $A$
- Goals  $G$

MDPs:

- Set of states  $S$
- Initial state  $I$
- Transition probabilities:
  - o  $P_a(s' | s)$
- Reward function  $r(s, a, s')$  in real
- Discount factor (gamma)  $\gamma \in [0, 1]$

Discounted reward:

$$V = r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4$$

$$= r_1 + \gamma(r_2 + \gamma r_3 + \gamma^2 r_4)$$

$$V_t = \frac{r_t + \gamma V_{t+1}}{\gamma(r_t + \gamma V_{t+1})}$$

$1.0 \times 0.9$

Probabilistic PDDL:

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

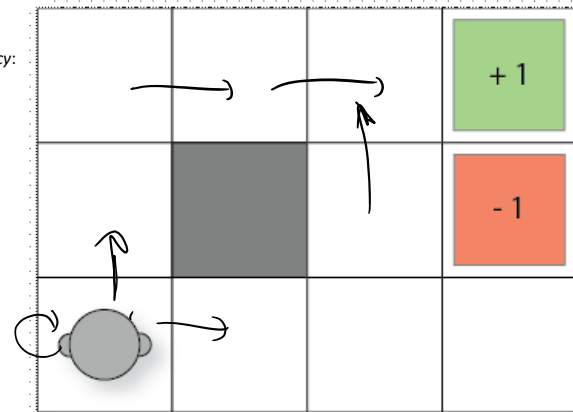
  (:action dunk-package
    :parameters (?pkg)
    :effect (and (when (bomb-in-package ?pkg)
                  (bomb-defused))
                (probabilistic 0.05 (toilet-clogged))))
```

Solution for MDP is a policy:

- at(0,0) => move\_right
- at(0,1) => move\_right
- at(0,2) => move\_right
- at(0,3) => stay
- at(1,0) => move\_up
- at(1,2) => move\_up
- at(1,3) => move\_up
- at(2,0) => move\_up
- at(2,1) => move\_left
- at(2,2) => move\_up
- at(2,3) => move\_left

$$\rightarrow \pi(s) \rightarrow a$$

policy



Expected return exercise:

You can steal:

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

Which do you steal? <https://pollev.com/timothymille936>

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

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

Expected return of  $a$

reward  $\rightarrow r'$

Bellman equations:

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

$$\underline{V(s)} = \max_{a \in A} \left( \sum_{s' \in S} \underbrace{P_{a,s|s}}_{\substack{\uparrow \\ \text{prob}}} (r(s,a,s') + \gamma \underline{V(s')}) \right)$$

$\uparrow$ 
reward  
function