

# Markov Decision Process MDP

- Create a policy when acting under uncertainty
- 1-Markov assumption: next state only depends on current state

## Components:

- Set of states:  $S$  Set of actions:  $A$
- Initial state:  $s_0$
- Transition model:  $T(s, a, s') \rightarrow [0, 1]$  if in  $s$  and perform  $a$ , what prob of getting to  $s'$   $P(s'/s, a)$
- Reward function:  $R(s) \rightarrow \text{real \#}$  How much value do you get for being in a state

Action VP

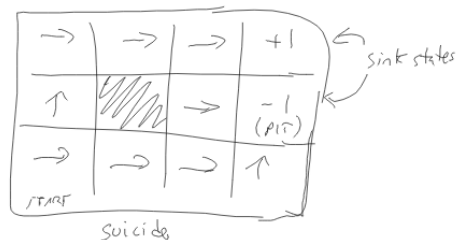
	(1,1)	(2,1)	(2,2)	...	(3,3)	...	(1,2)
(1,1)	0.1	0.1			(0.0)		0.8
(2,1)	0.1	0.8			0.1		0.0
(2,2)							
...							
(3,3)							
...							
(1,2)							

## MDP

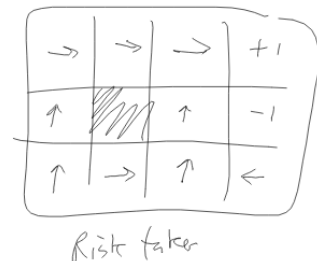
$R(s) \rightarrow R$   
Reward defines what optimal is.

- Solution is a policy  $\pi: s \rightarrow a$
- $\pi(s)$  is the action to take from state  $s$
- $\pi^*(s)$  is the optimal action to take from state  $s$
- Maximizes reward over time, giving the highest expected utility

$$R(s) \leq -1.7$$



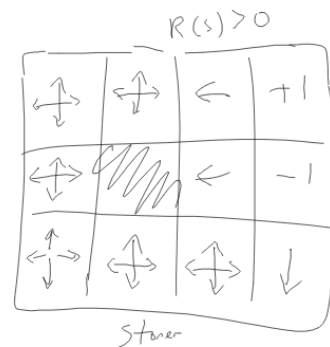
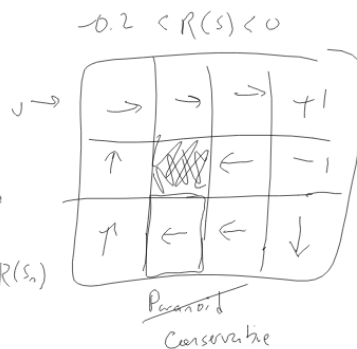
$$-0.43 < R(s) < -0.06$$



## Calculating utility

- What is it worth to the agent to be in a state?
- Utility of a state  $s$  is the amount of reward accumulated since starting in state  $s_0$  before arriving in  $s$
- Depends on how you got there
- Utility of a plan acts like a heuristic, indicating which route is better

$$U([s_0, s_1, s_2, \dots, s_n]) = R(s_0) + R(s_1) + \dots + R(s_n)$$



$R(s, a)$

## Calculating utility

### Discounted reward:

- $U_n([s_0, s_1, \dots, s_n]) = \underline{R(s_0)} + \underline{\gamma R(s_1)} + \underline{\gamma^2 R(s_2)} + \dots + \underline{\gamma^n R(s_n)}$
- $\gamma$  is discount factor,  $[0, 1]$

- If you knew the utility of every state, the policy is just the action that moves the agent to the successor with the highest utility

0.812	0.868	0.918	+1
0.762		0.660	-1
0.705 <i>start</i>	0.655	0.611	0.388

$U(\text{all future sequences})$

$U(R)$

$U(U)$

$U(UR)$

$U(DUR)$

Max

$$\pi^*(s) = \underset{a \in A}{\operatorname{argmax}} \left( \sum_{s' \in S} (T(s,a,s') \times U(s')) \right)$$

for all other states

Probability of getting to that state if a is executed

Utility of Successor

(?)

- Need to compute  $U(s)$
- Utility of state  $s$  is the immediate reward  $R(s)$  plus expected discounted utility of future states

- **Bellman Equation:**

$$U(s) = R(s) + \gamma \times \max_{a \in A} \left( \sum_{s' \in S} T(s, a, s') \times U(s') \right)$$

The equation is annotated with handwritten notes and arrows:
 

- An arrow points from the handwritten word "me" to the underlined  $U(s)$ .
- An arrow points from the handwritten word "my reward" to  $R(s)$ .
- An arrow points from the handwritten word "discount" to  $\gamma$ .
- An arrow points from the handwritten phrase "probability of successor state" to  $T(s, a, s')$ .
- An arrow points from the handwritten word "Recursive!" to the circled  $U(s')$ .
- A separate circle containing a question mark "?" is connected by a line to the circled  $U(s')$ .

# Value Iteration

- Start with arbitrary utility values and iteratively update the values until they converge on true utility
- Recall: utility of a state is a function on immediate reward and the expected utilities of all neighbors
- $U_i(s)$  is the utility of state  $s$  at the  $i$ th iteration

- **Bellman update:**

$$U_{i+1}(s) = R(s) + \gamma \times \max_{a \in A} \left( \sum_{s' \in S} (T(s, a, s') \times U_i(s')) \right)$$

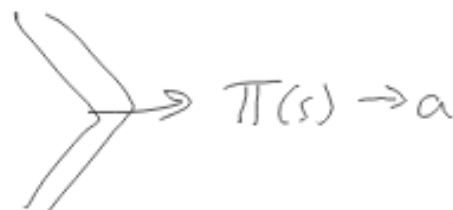
$\uparrow$  my utility at  $i+1$        $\uparrow$  my reward       $\uparrow$  Utility of my successors at the previous iteration

MDP

States  $S$   
 Actions  $A$   
 initial state  $s_0$

$T(s', a, s)$

$R(s)$  ( $R(s, a)$ )



$\pi(s) \rightarrow a$

## Value Iteration Algorithm

- Set  $U_0(s)$  to arbitrary starting value for all  $s$
- Do until utilities converge (difference less than  $\epsilon$ ):
  - $i = i + 1$
  - For each state  $s$  do:
    - Compute  $U_i$  using the Bellman update (using  $U_{i-1}$ ) 