

Probabilistic Robotics

The Markov Decision Process

**(adapted from Thrun, Burgard,
Fox & Russell/Norvig Chap 17)**

Where are We ?

- Estimator
 - Input: sensor readings, commanded actions
 - Output: Belief distribution over all possible values of the robot's state

Sensing model
Action model
Bayes filter

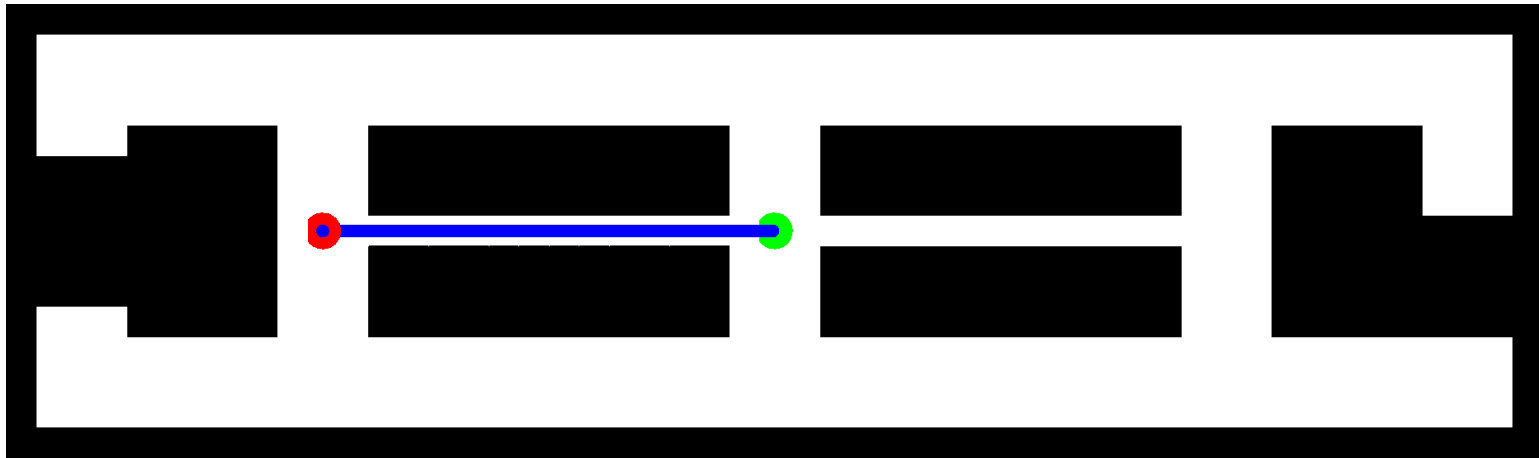
- Controller/Planner
 - Input: Robot's current state, goal state
 - Output: action to take next

TODAY

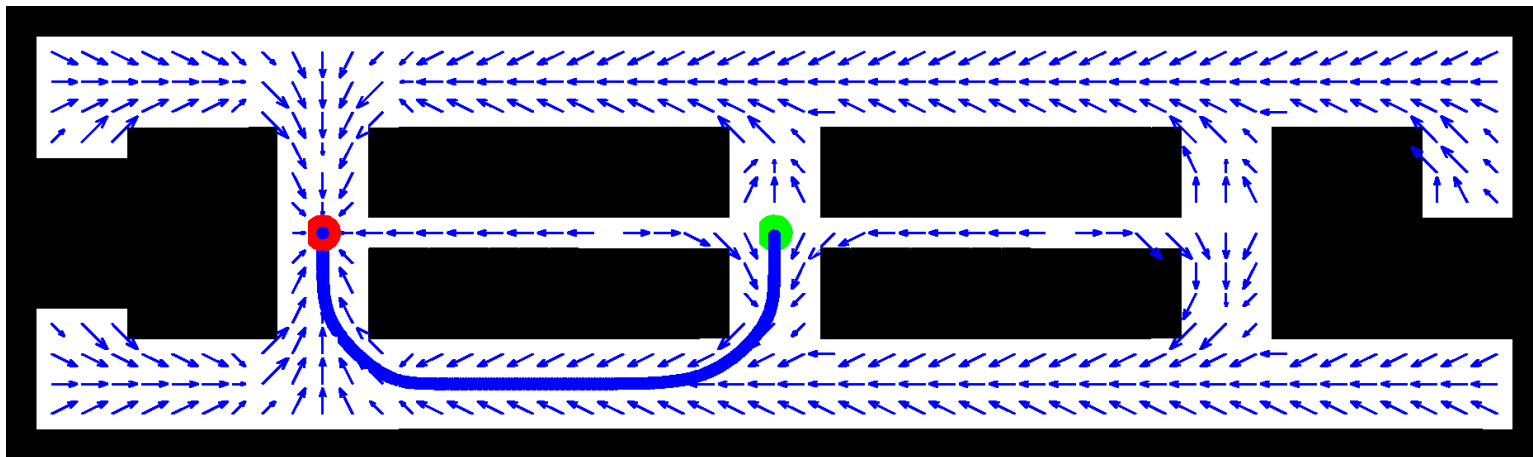
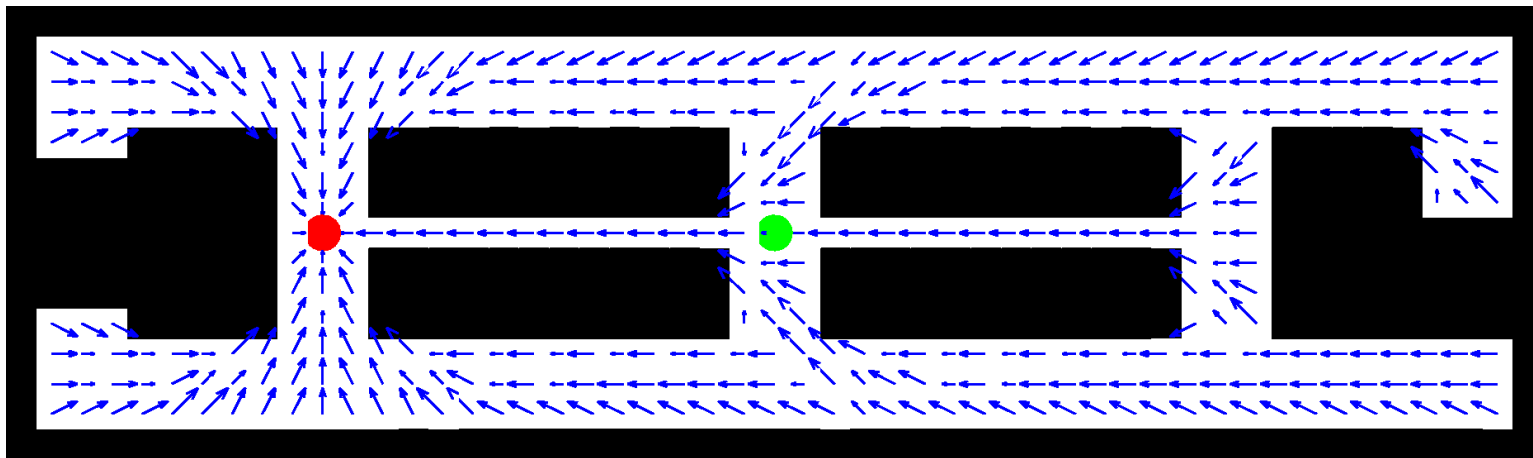
Problem Classes

- Controller/Planner
 - Input: Robot's current state, goal state
 - Output: action to take next
- Deterministic vs. stochastic actions
- Full vs. partial observability

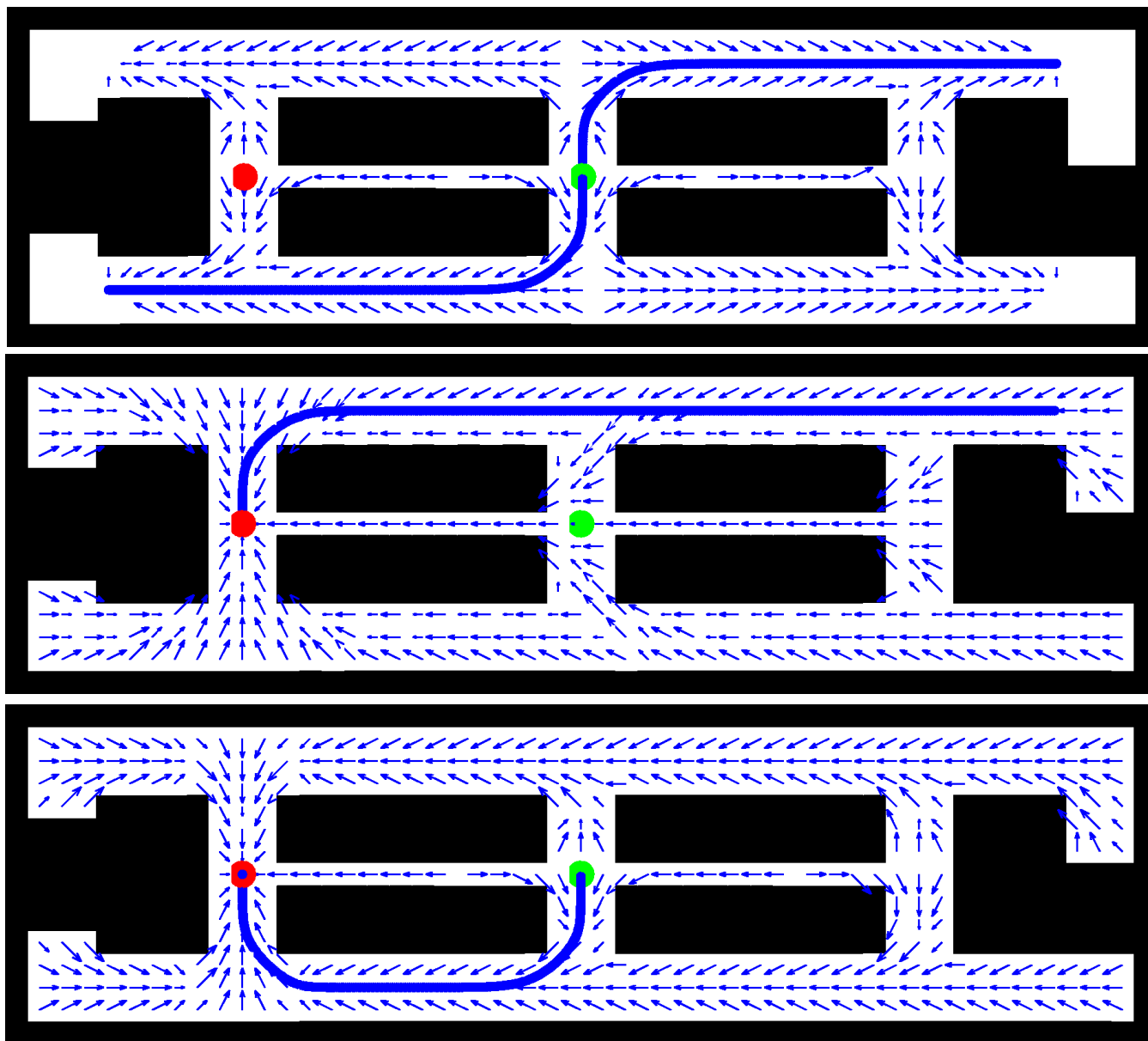
Deterministic, fully observable



Stochastic, Fully Observable



Stochastic, Partially Observable



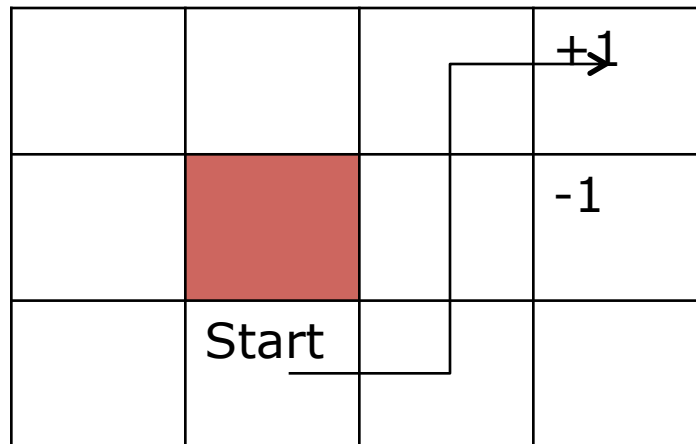
Russell/Norvig (Chap 17)

- Robot Motion planning example
 - End up at a goal (+1)
 - Avoid some place(s) (-1)
 - Moving has a cost (e.g., -0.04 per cell)

			+1
			-1
	Start		

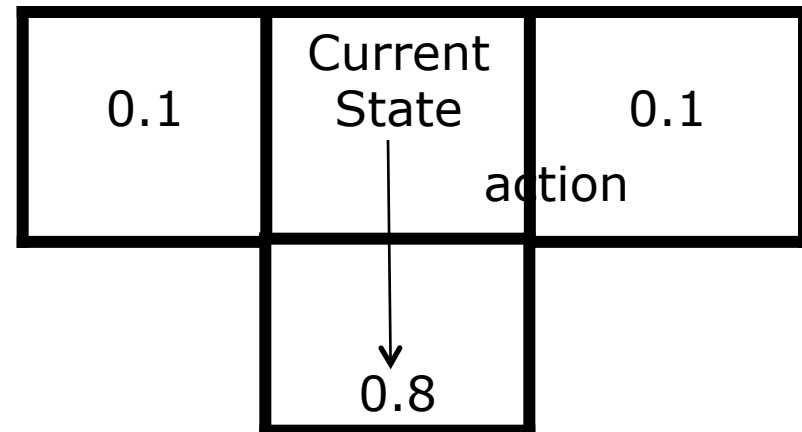
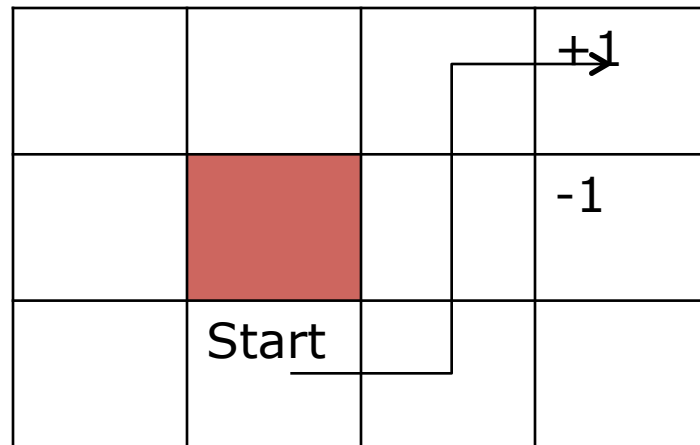
Deterministic case

- Just take the shortest path (e.g. Dijkstra)



Non-Deterministic Case

- Stochastic Action Model
- What is the best way to go from start to goal ?



Pick the highest 'utility' path

- Longer
- Smaller probability of ending up in bad state
- Deterministic case is special case of non-deterministic with utility = $(1/\text{distance to goal})$



Utility and Policy

- Compute for every state: how useful is each state for the overall task/problem ?

$$U(i) = R(i) + \max_a \sum_j M_{ij}^a \cdot U(j)$$

- Policy: Mapping from states to actions (for all states)

Markov Decision Process/ Problem

- Compute optimal policy in a stochastic setting with a known transition model
- Markov property: transition probabilities depend only on current state

Optimal Policy

- If we know utility we can compute optimal policy
- How to compute utilities ?

$$policy^*(i) = \operatorname{argmax}_a \sum_j M_{ij}^a \cdot U(j)$$

M_{ij}^a = Probability of reaching state j from state i with action a .

$U(j)$ = Utility of state j .

Compute Utility Iteratively

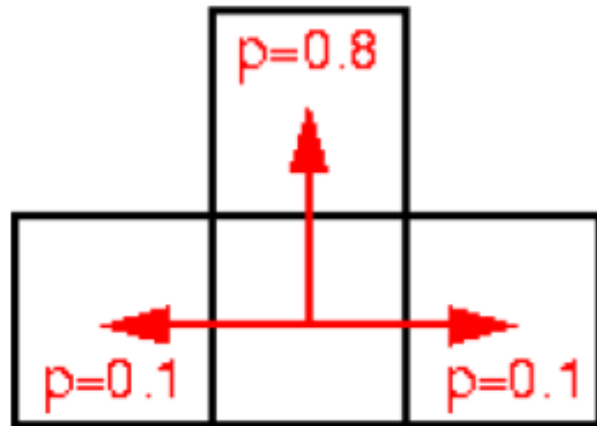
- Can go to infinity, in practice stop when some threshold reached

$$U_{t+1}(i) = R(i) + \max_a \sum_j M_{ij}^a \cdot U_t(j)$$

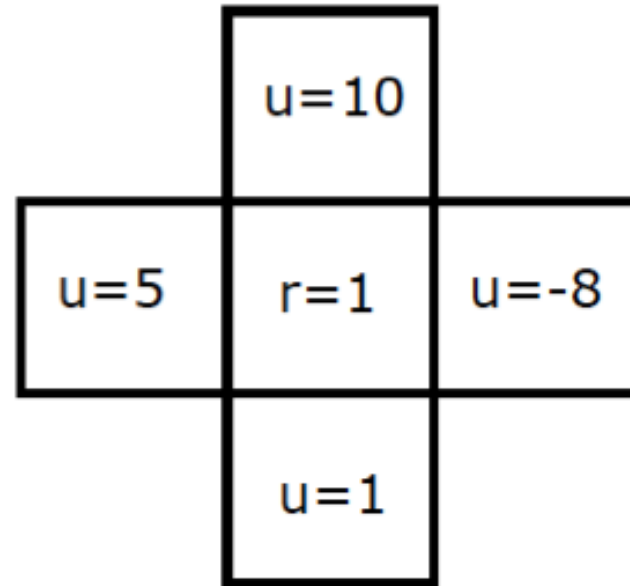
Value Iteration Example

$$U_{t+1}(i) = R(i) + \max_a \sum_j M_{ij}^a \cdot U_t(j)$$

(desired action = North)

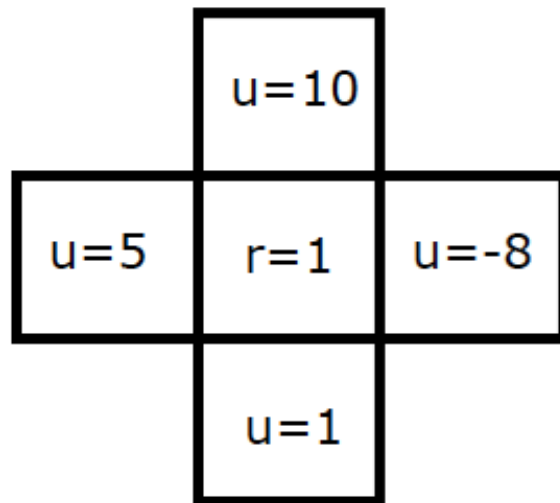


Transition Model



State Space
(u =utility, r =reward)

Value Iteration Example



$$\begin{aligned} U_{t+1}(i) &= R(i) + \max_a \sum_j M_{ij}^a \cdot U_t(j) \\ &= \text{reward} + \max\{ \\ &\quad 0.1 \cdot 1 + 0.8 \cdot 5 + 0.1 \cdot 10 \quad (\leftarrow), \\ &\quad 0.1 \cdot 5 + 0.8 \cdot 10 + 0.1 \cdot -8 \quad (\uparrow), \\ &\quad 0.1 \cdot 10 + 0.8 \cdot -8 + 0.1 \cdot 1 \quad (\rightarrow), \\ &\quad 0.1 \cdot -8 + 0.8 \cdot 1 + 0.1 \cdot 5 \quad (\downarrow)\} \\ &= 1 + \max\{5.1 (\leftarrow), 7.7 (\uparrow), \\ &\quad -5.3 (\rightarrow), 0.5 (\downarrow)\} \\ &= 1 + 7.7 \\ &= 8.7 \end{aligned}$$

Finding the Optimal Policy

$$policy^*(i) = \operatorname{argmax}_a \sum_j M_{ij}^a \cdot U^*(j)$$

			+1
			-1

1. The given environment.

0.812	0.868	0.912	+1
0.762		0.660	-1
0.705	0.655	0.611	0.388

2. Calculate Utilities.

Why does this arrow point left (not up?)

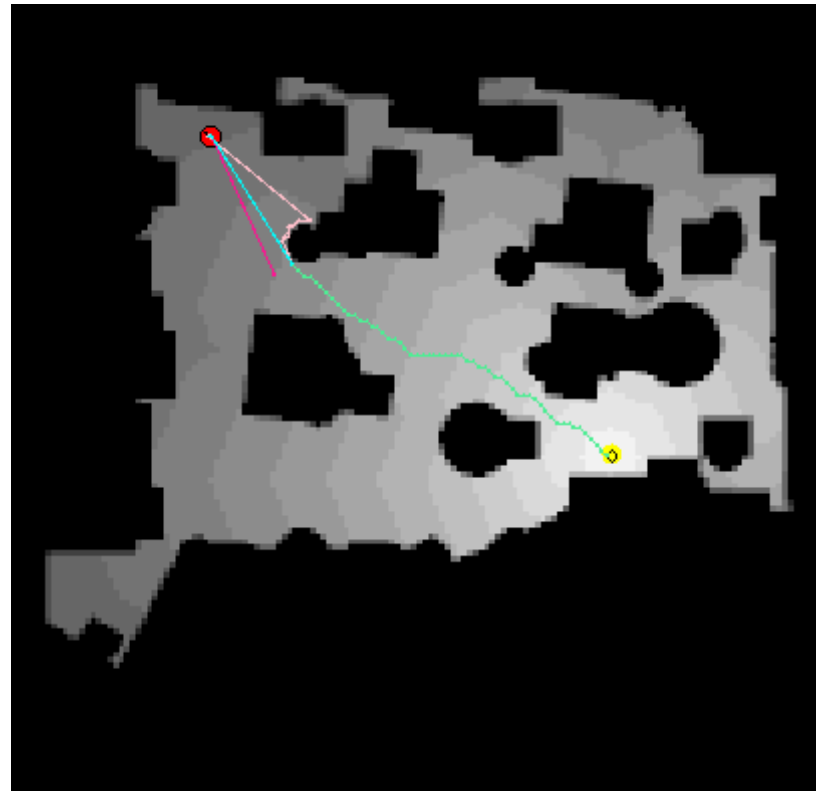
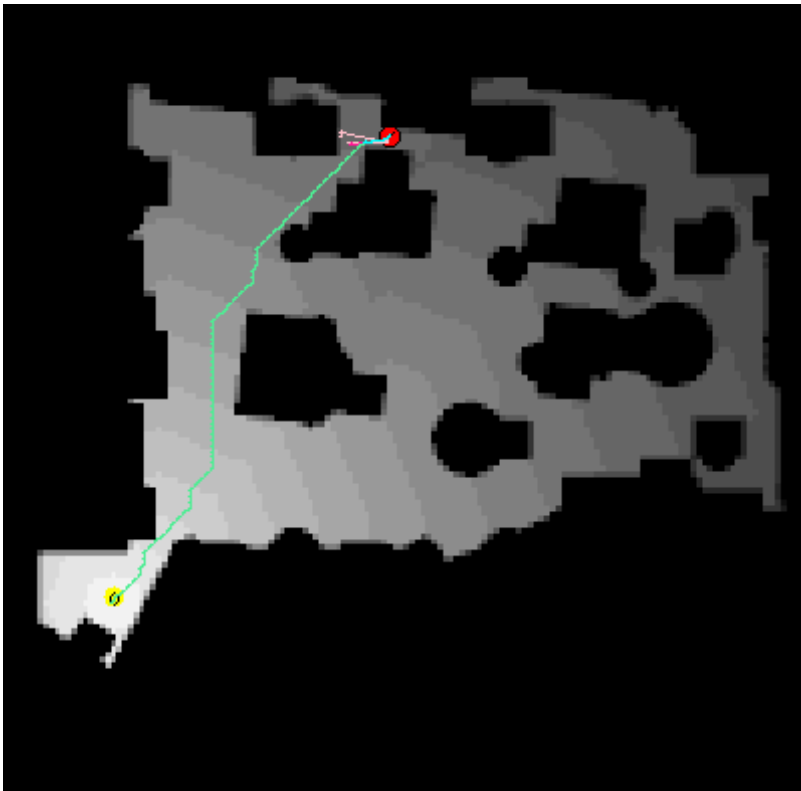
→	→	→	+1
↑		↑	-1
↑	←	←	←

3. Extract optimal policy.

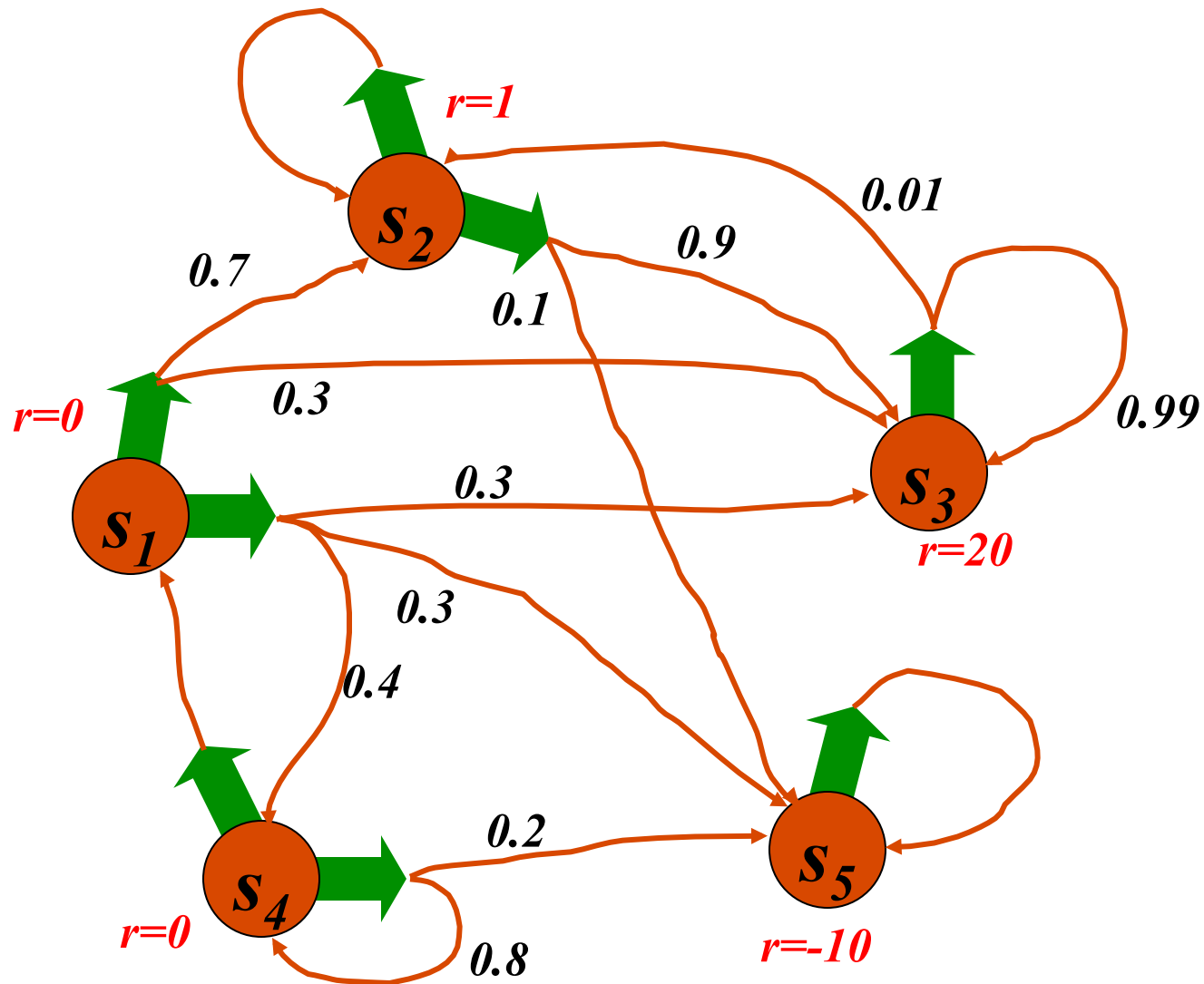
			+1
			-1

4. Execute actions.

Value Iteration for Motion Planning



Markov Decision Process (MDP)



Markov Decision Process (MDP)

- **Given:**
- States x
- Actions u
- Transition probabilities $p(x' | u, x)$
- Reward / payoff function $r(x, u)$

- **Wanted:**
- Policy $\pi(x)$ that maximizes the future expected reward

Rewards and Policies

- Policy (general case):

$$\pi: z_{1:t-1}, u_{1:t-1} \rightarrow u_t$$

- Policy (fully observable case):

$$\pi: x_t \rightarrow u_t$$

- Expected cumulative payoff:

$$R_T = E \left[\sum_{\tau=1}^T \gamma^\tau r_{t+\tau} \right]$$

- $T=1$: greedy policy
- $T>1$: finite horizon case, typically no discount
- $T=\infty$: infinite-horizon case, finite reward if discount < 1

Policies contd.

- Expected cumulative payoff of policy:

$$R_T^\pi(x_t) = E \left[\sum_{\tau=1}^T \gamma^\tau r_{t+\tau} \mid u_{t+\tau} = \pi(z_{1:t+\tau-1}, u_{1:t+\tau-1}) \right]$$

- Optimal policy:

$$\pi^* = \operatorname{argmax}_{\pi} R_T^\pi(x_t)$$

- 1-step optimal policy:

$$\pi_1(x) = \operatorname{argmax}_u r(x, u)$$

- Value function of 1-step optimal policy:

$$V_1(x) = \gamma \max_u r(x, u)$$

2-step Policies

- Optimal policy:

$$\pi_2(x) = \operatorname{argmax}_u \left[r(x, u) + \int V_1(x') p(x' | u, x) dx' \right]$$

- Value function:

$$V_2(x) = \gamma \max_u \left[r(x, u) + \int V_1(x') p(x' | u, x) dx' \right]$$

T-step Policies

- Optimal policy:

$$\pi_T(x) = \operatorname{argmax}_u \left[r(x, u) + \int V_{T-1}(x') p(x' | u, x) dx' \right]$$

- Value function:

$$V_T(x) = \gamma \max_u \left[r(x, u) + \int V_{T-1}(x') p(x' | u, x) dx' \right]$$

Infinite Horizon

- Optimal policy:

$$V_{\infty}(x) = \gamma \max_u \left[r(x, u) + \int V_{\infty}(x') p(x' | u, x) dx' \right]$$

- Bellman equation
- Fix point is optimal policy
- Necessary and sufficient condition

Value Iteration

- for all x do

$$\hat{V}(x) \leftarrow r_{\min}$$

- endfor

- repeat until convergence

- for all x do

$$\hat{V}(x) \leftarrow \gamma \max_u \left[r(x, u) + \int \hat{V}(x') p(x' | u, x) dx' \right]$$

- endfor

- endrepeat

$$\pi(x) = \operatorname{argmax}_u \left[r(x, u) + \int \hat{V}(x') p(x' | u, x) dx' \right]$$

Value Function and Policy Iteration

- Often the optimal policy has been reached long before the value function has converged.
- Policy iteration calculates a new policy based on the current value function and then calculates a new value function based on this policy.
- This process often converges faster to the optimal policy.

Where are We ?

- Estimator

- Input: sensor readings, commanded actions
- Output: Belief distribution over all possible values of the robot's state

Sensing model
Action model
Bayes filter

- Controller/Planner

- Input: Robot's current state, goal state
- Output: action to take next

Action model
MDP/VI

Next time: POMDPs

- In POMDPs we apply the very same idea as in MDPs.
- **Since the state is not observable**, the agent has to **make its decisions based on** the belief state which is a **posterior distribution over states**.
- Let b be the belief of the agent about the state under consideration.
- POMDPs compute a **value function over belief space**:

$$V_T(b) = \gamma \max_u \left[r(b, u) + \int V_{T-1}(b') p(b' | u, b) db' \right]$$