# CSE-571
# Robotics

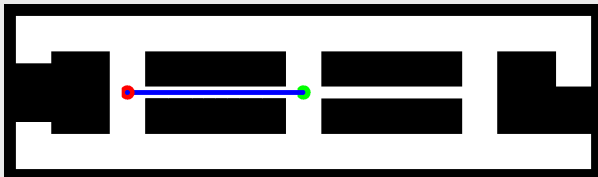**Planning and Control:**

**Markov Decision Processes**

## Problem Classes

- Deterministic vs. stochastic actions

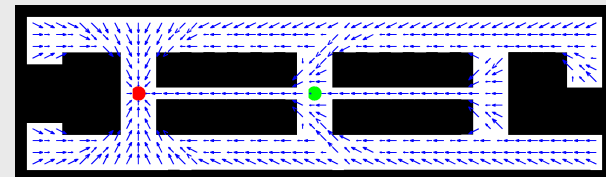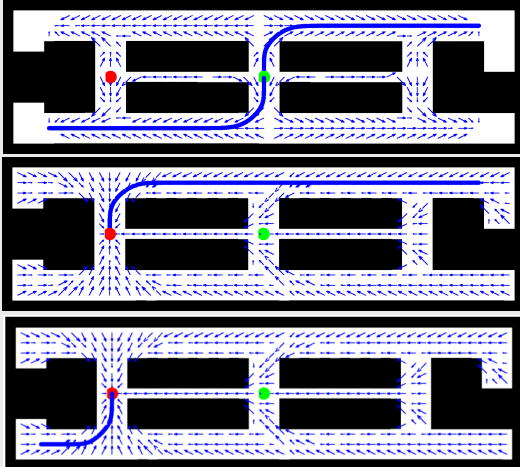- Full vs. partial observability

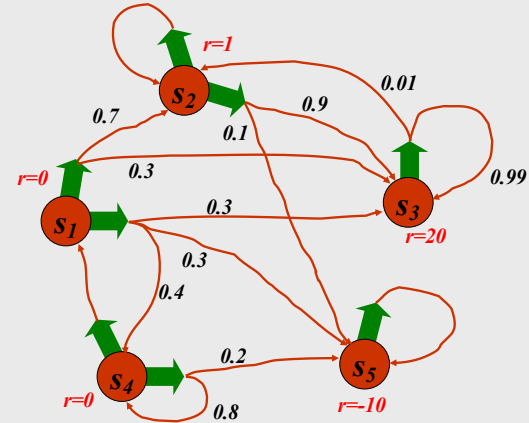## Deterministic, fully observable

## Stochastic, Fully Observable

## Stochastic, Partially Observable



5

## Markov Decision Process (MDP)



6

## 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

7

## Rewards and Policies

- Policy (general case):

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

- Policy (fully observable case):

$$\pi: \quad x_t \to u_t$$

- 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

8

## 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:

$$\pi^* = \underset{\pi}{\mathrm{argmax}} \quad R_T^\pi(x_t)$$

- 1-step optimal policy:

$$\pi_1(x) = \underset{u}{\mathrm{argmax}} \quad r(x,u)$$

- Value function of 1-step optimal policy:

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

## 2-step Policies

- Optimal policy:

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

- 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:

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

- 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

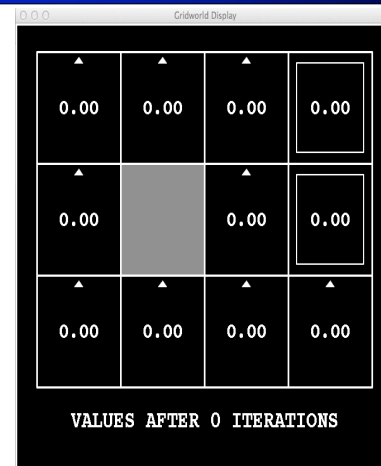$$\hat{V}(x) \leftarrow \gamma \max_{u} \left[ r(x,u) + \int \hat{V}(x')p(x'|u,x)dx' \right]$$

  - endfor
- endrepeat

$$\pi(x) = \underset{u}{\operatorname{argmax}} \left[ r(x,u) + \int \hat{V}(x')p(x'|u,x)dx' \right]$$
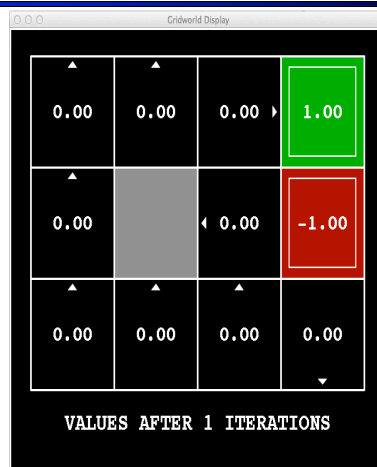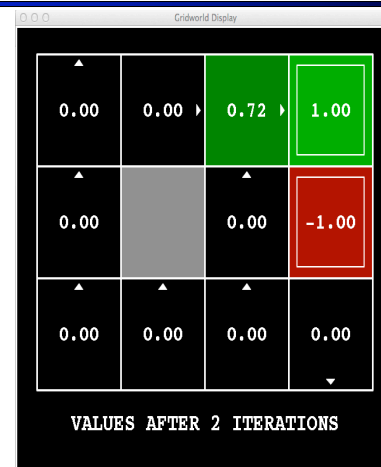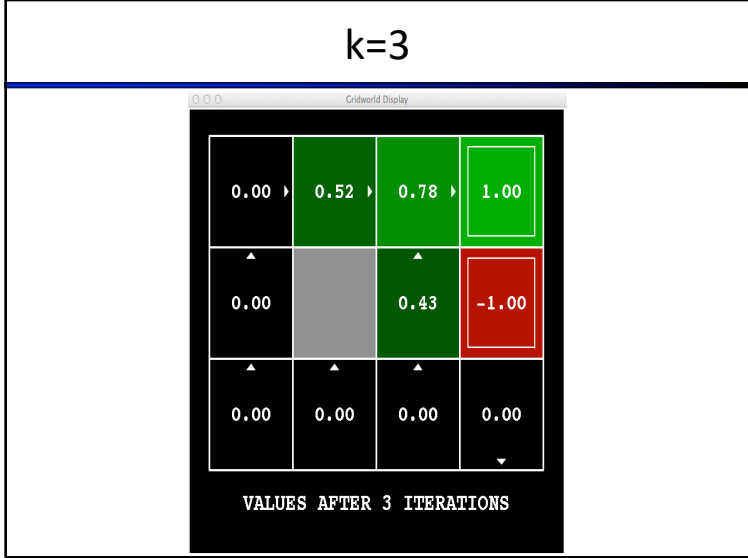
13

---

## k=0



Noise = 0.2
Discount = 0.9
Living reward = 0

14

---

## k=1



Noise = 0.2
Discount = 0.9
Living reward = 0

15

---

## k=2



16

## k=3

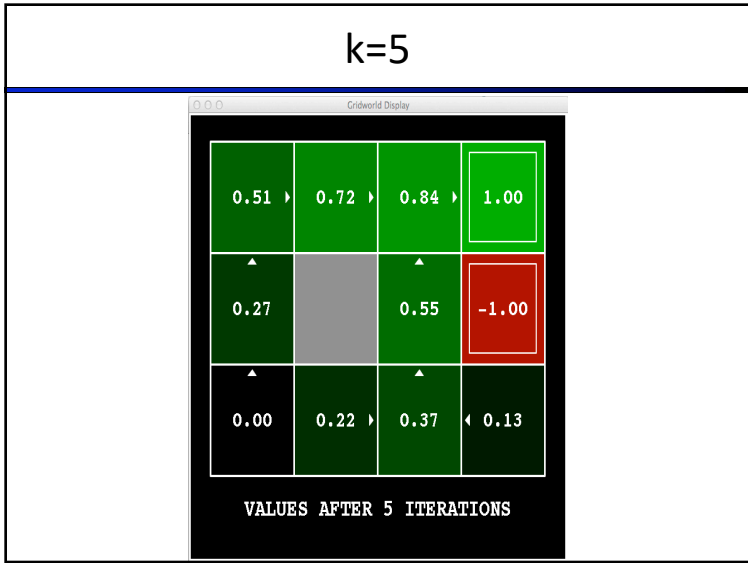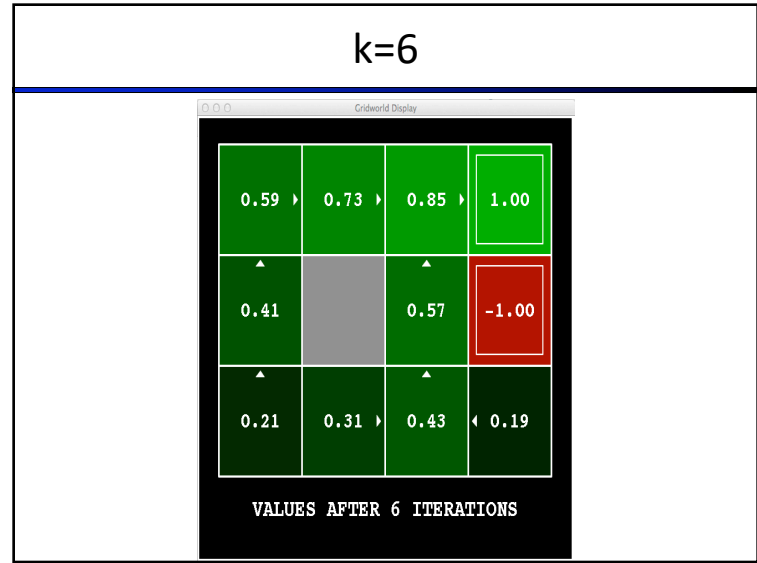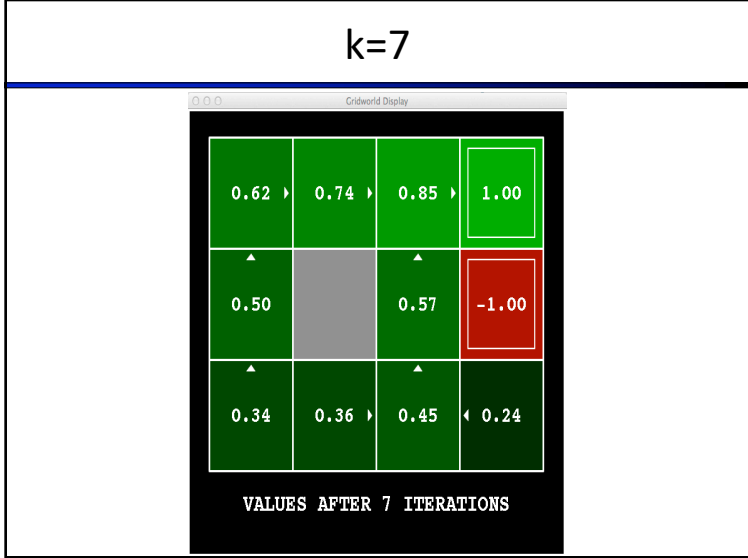

| 0.00 ▸ | 0.52 ▸ | 0.78 ▸ | 1.00 |
|---|---|---|---|
| ▲ 0.00 | | ▲ 0.43 | -1.00 |
| ▲ 0.00 | ▲ 0.00 | ▲ 0.00 | 0.00 ▾ |

VALUES AFTER 3 ITERATIONS

17

## k=4



| 0.37 ▸ | 0.66 ▸ | 0.83 ▸ | 1.00 |
|---|---|---|---|
| ▲ 0.00 | | ▲ 0.51 | -1.00 |
| ▲ 0.00 | 0.00 ▸ | ▲ 0.31 | ◂ 0.00 |

VALUES AFTER 4 ITERATIONS

18

## k=5



| 0.51 ▸ | 0.72 ▸ | 0.84 ▸ | 1.00 |
|---|---|---|---|
| ▲ 0.27 | | ▲ 0.55 | -1.00 |
| ▲ 0.00 | 0.22 ▸ | ▲ 0.37 | ◂ 0.13 |

VALUES AFTER 5 ITERATIONS

19

## k=6



| 0.59 ▸ | 0.73 ▸ | 0.85 ▸ | 1.00 |
|---|---|---|---|
| ▲ 0.41 | | ▲ 0.57 | -1.00 |
| ▲ 0.21 | 0.31 ▸ | ▲ 0.43 | ◂ 0.19 |

VALUES AFTER 6 ITERATIONS

20

5

## k=7



VALUES AFTER 7 ITERATIONS

## k=8



VALUES AFTER 8 ITERATIONS

## k=9



VALUES AFTER 9 ITERATIONS

## k=10



VALUES AFTER 10 ITERATIONS

21

22

23

24

**k=11**

Gridworld Display

| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| ▲ 0.56 | | ▲ 0.57 | −1.00 |
| ▲ 0.48 | ◀ 0.42 | ▲ 0.47 | ◀ 0.27 |

VALUES AFTER 11 ITERATIONS

25



**k=12**

Gridworld Display

| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| ▲ 0.57 | | ▲ 0.57 | −1.00 |
| ▲ 0.49 | ◀ 0.42 | ▲ 0.47 | ◀ 0.28 |

VALUES AFTER 12 ITERATIONS

26



**k=100**

Gridworld Display

| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| ▲ 0.57 | | ▲ 0.57 | −1.00 |
| ▲ 0.49 | ◀ 0.43 | ▲ 0.48 | ◀ 0.28 |

VALUES AFTER 100 ITERATIONS
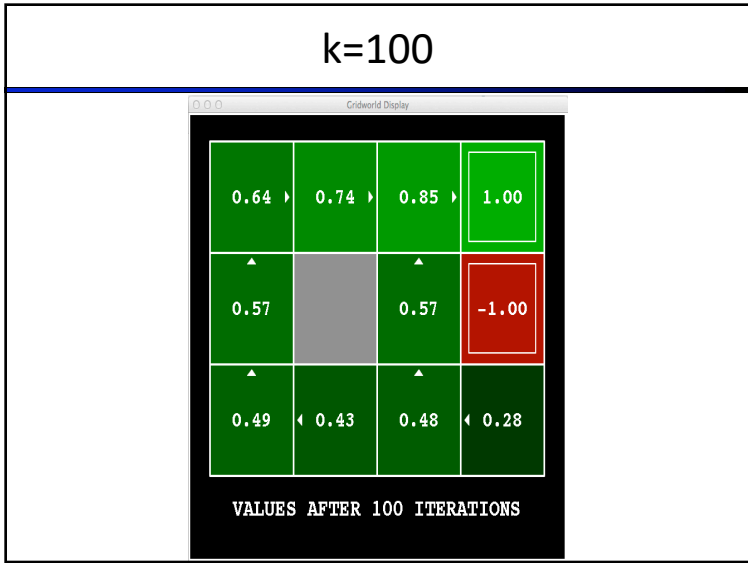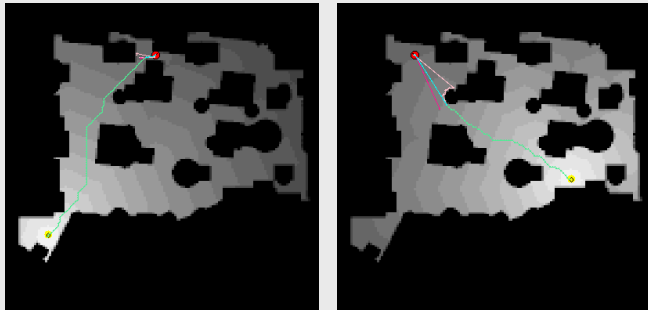
27

## Value Function and Policy

- Each step takes O(|A| |S| |S|) time.
- Number of iterations required is polynomial in |S|, |A|, 1/(1-gamma)



28

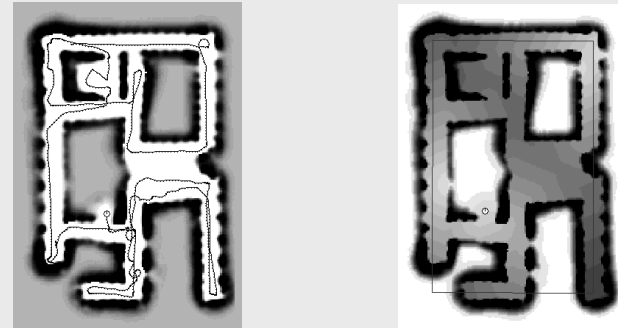7

## Value Iteration for Motion Planning
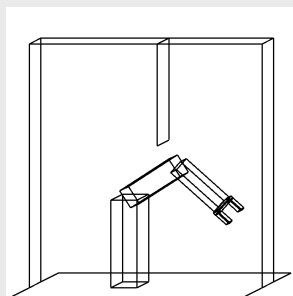**(assumes knowledge of robot's location)**

## Frontier-based Exploration

- Every unknown location is a target point.
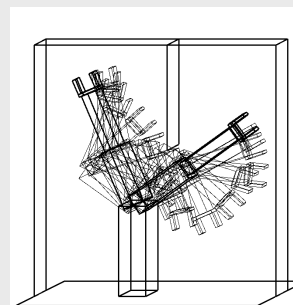
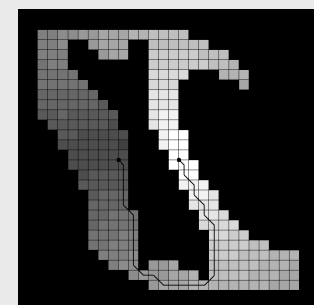## Manipulator Control



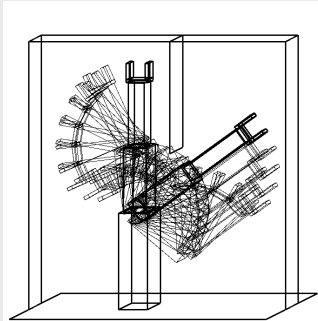Arm with two joints    Configuration space
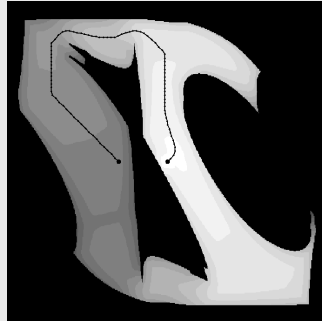
## Manipulator Control Path



State space    Configuration space

## Manipulator Control Path



State space      Configuration space

33

## 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.
- For finite horizon problems, the resulting value functions are piecewise linear and convex.
- In each iteration the **number of linear constraints grows exponentially**.
- Full fledged POMDPs have only been applied to very small state spaces with small numbers of possible observations and actions.
- **Approximate solutions are becoming more and more capable**.

34