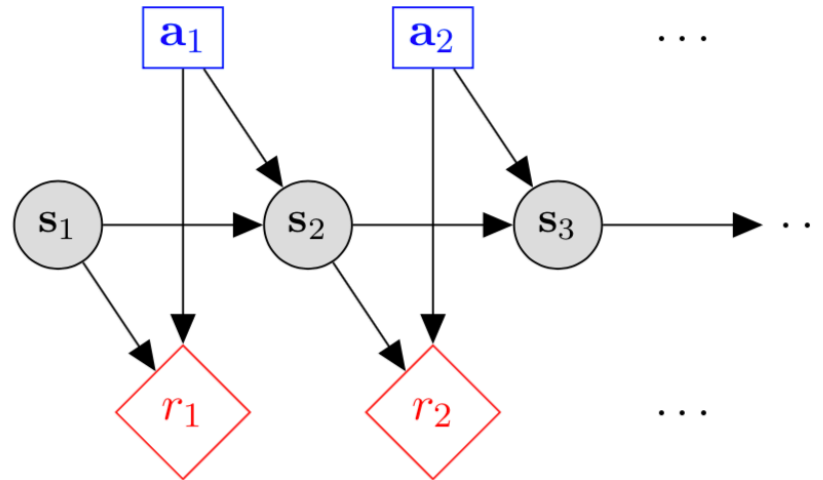# CS181: Introduction to Machine Learning
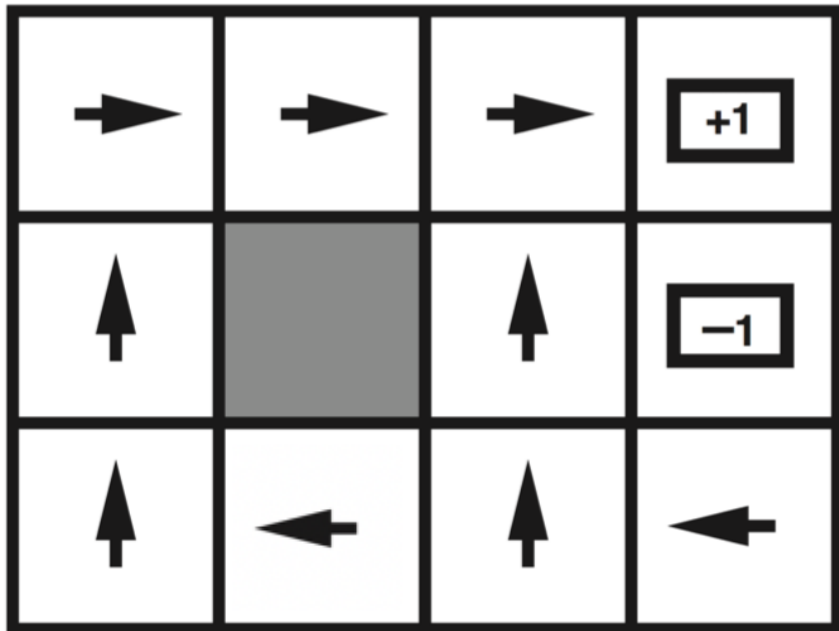
# Lecture 21 (MDPs and RL)

## Spring 2021

Finale Doshi-Velez and David C. Parkes
Harvard Computer Science

# Example: House cleaning robot



- States: physical location, objects in environment
- Actions: move, pick-up, drop, ...
- Reward: $+1$ if pick up dirty clothes, -1 if break dish, ...
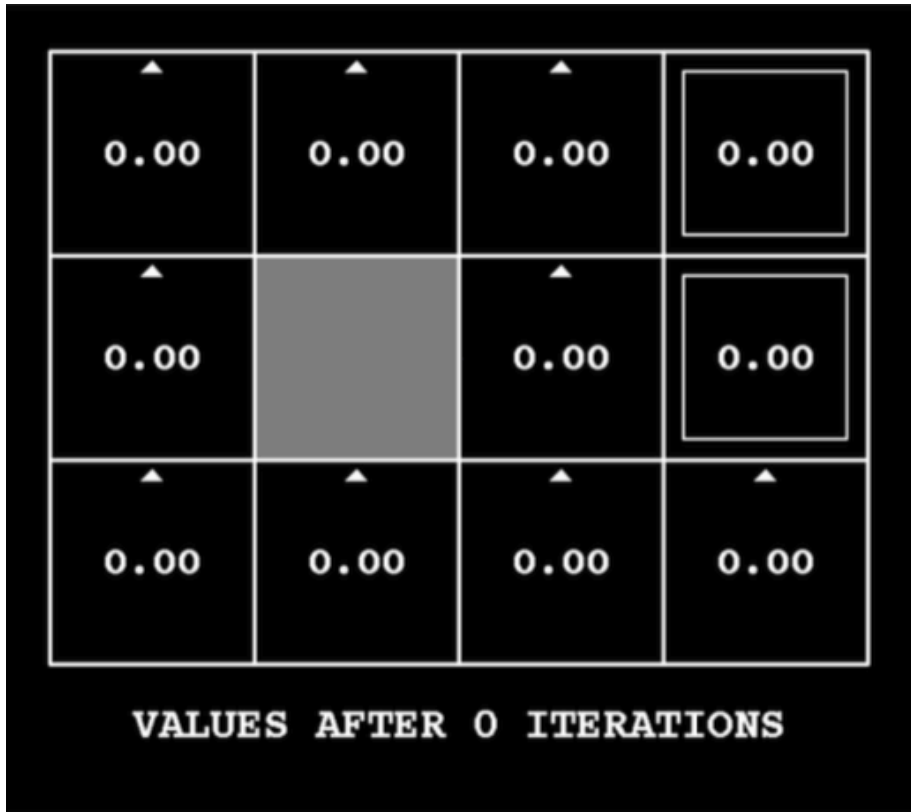- Transition model: describe actuators and uncertain environment
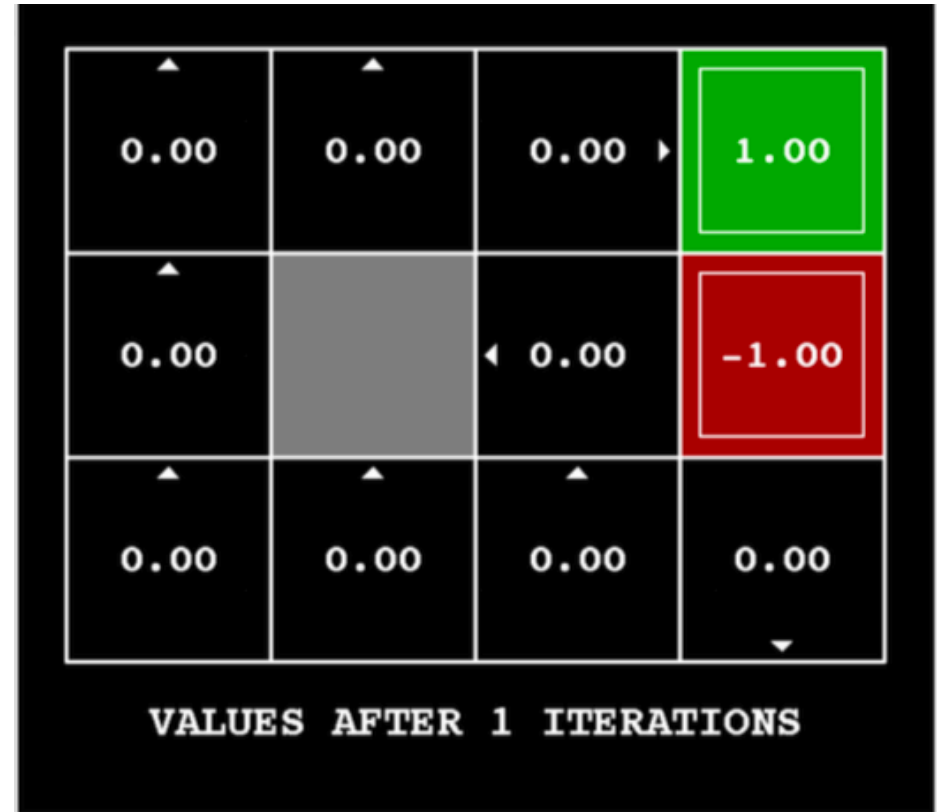
# GridWorld



(D. Klein, P. Abbeel)

- r(s,a) = 0, except states (1,4), (2,4). In these states get +1 or -1 when take ANY action. Then no more actions

- Bounce off obstacles. Actuator has 20% noise; e.g., w/ prob 0.1 goes L, prob 0.1 goes R when moving U

- Discounting 0.9 (r + 0.9 r + $0.9^2$ r + …)

# VI, GridWorld

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



VALUES AFTER 0 ITERATIONS



VALUES AFTER 1 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 1 ITERATIONS



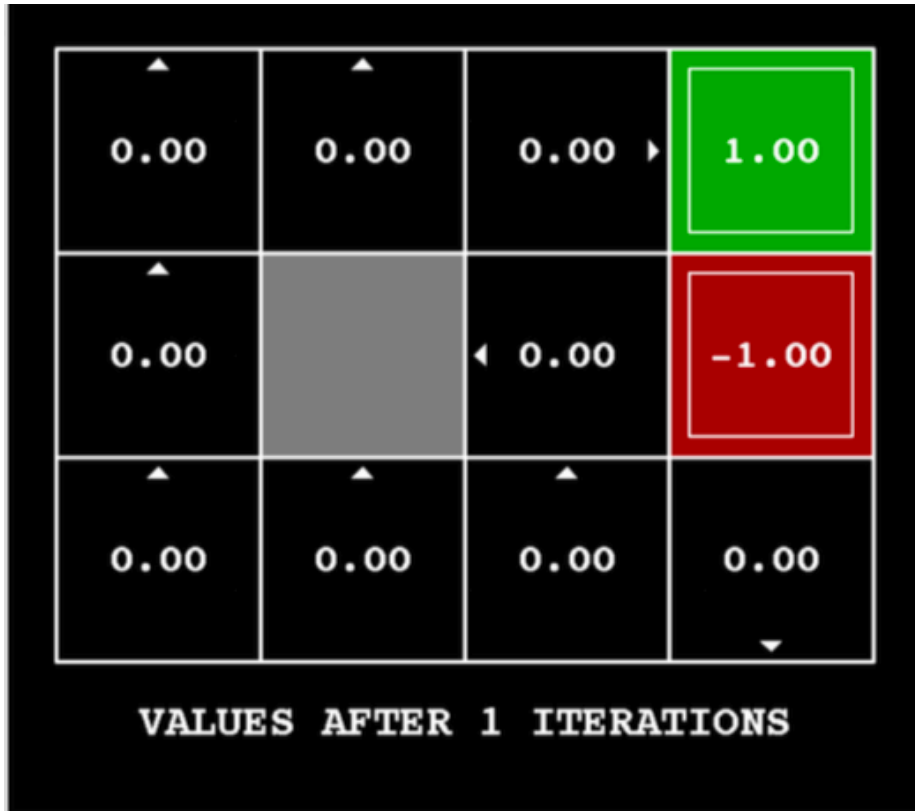VALUES AFTER 2 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



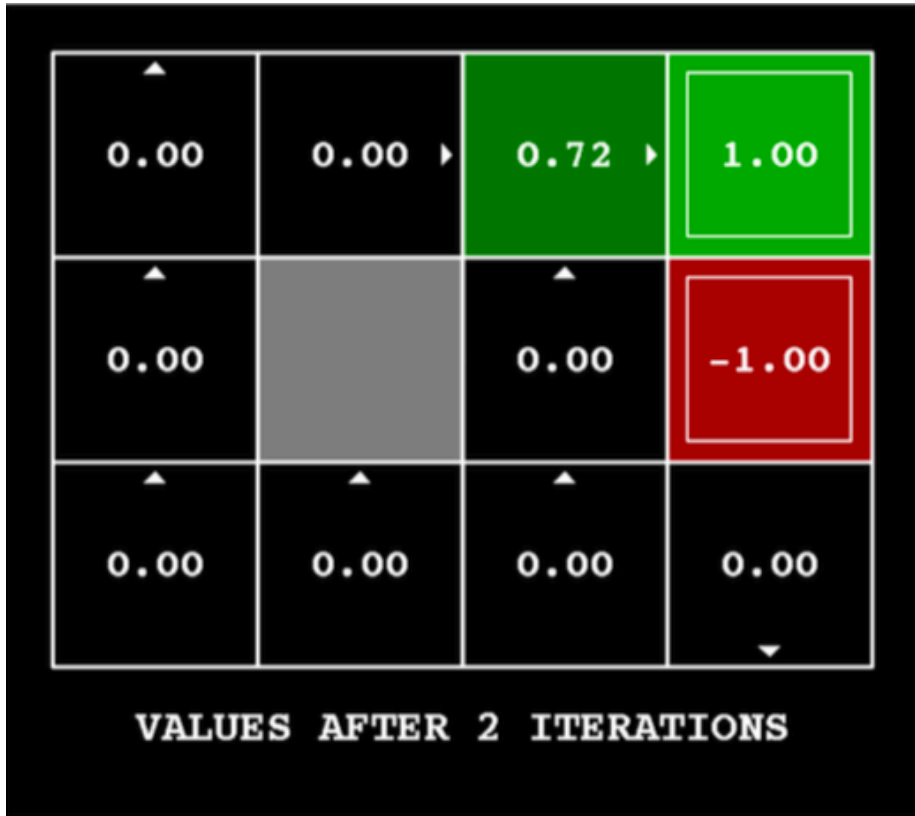VALUES AFTER 2 ITERATIONS



VALUES AFTER 3 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 3 ITERATIONS



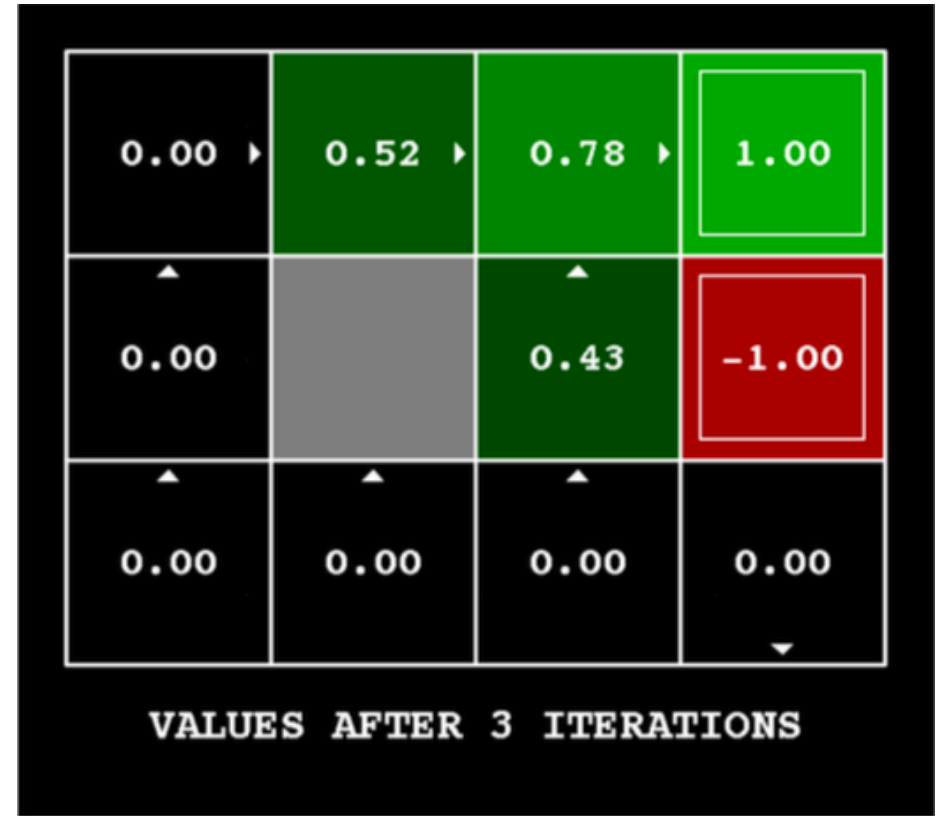VALUES AFTER 4 ITERATIONS

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 4 ITERATIONS



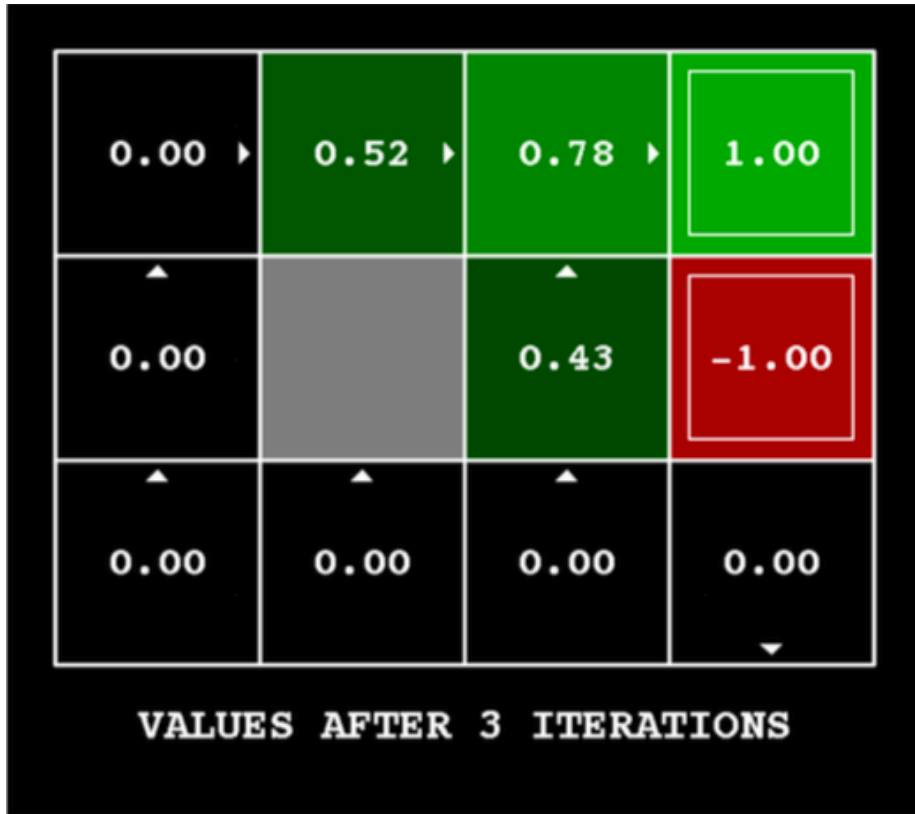VALUES AFTER 5 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 5 ITERATIONS



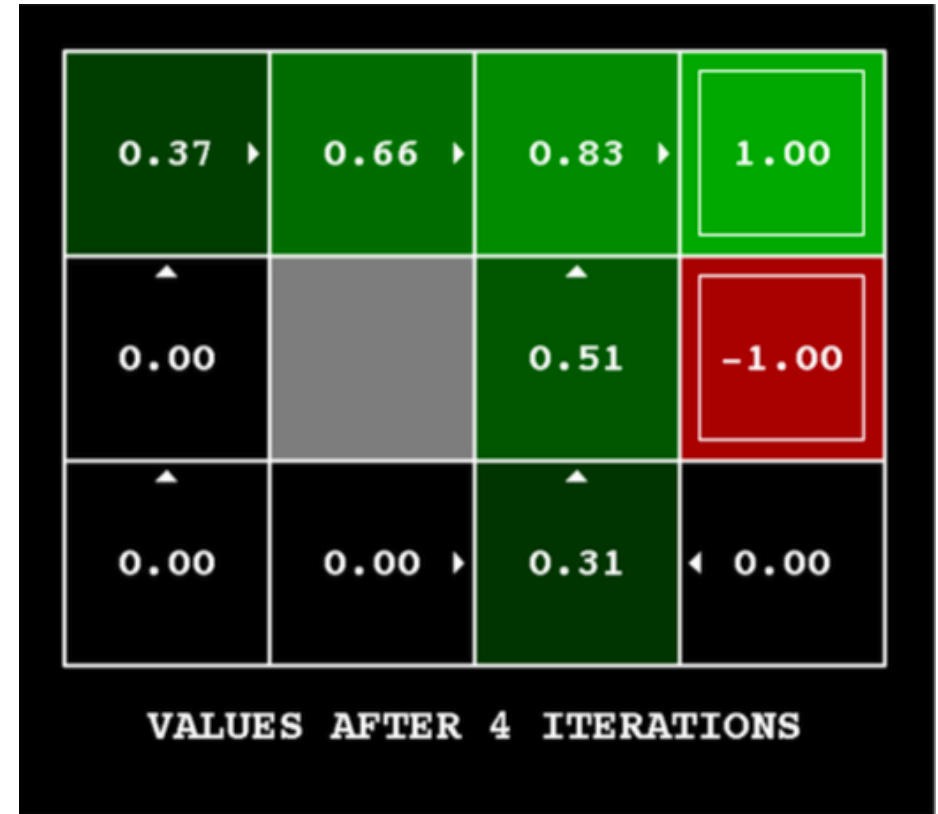VALUES AFTER 6 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 6 ITERATIONS



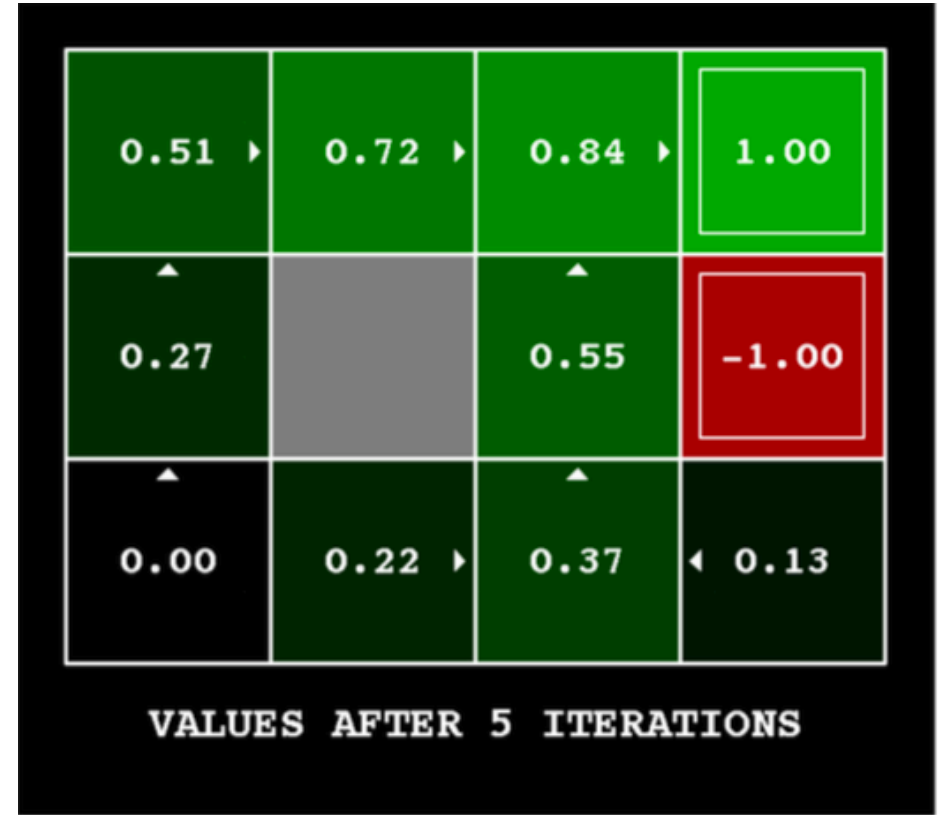VALUES AFTER 7 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 7 ITERATIONS

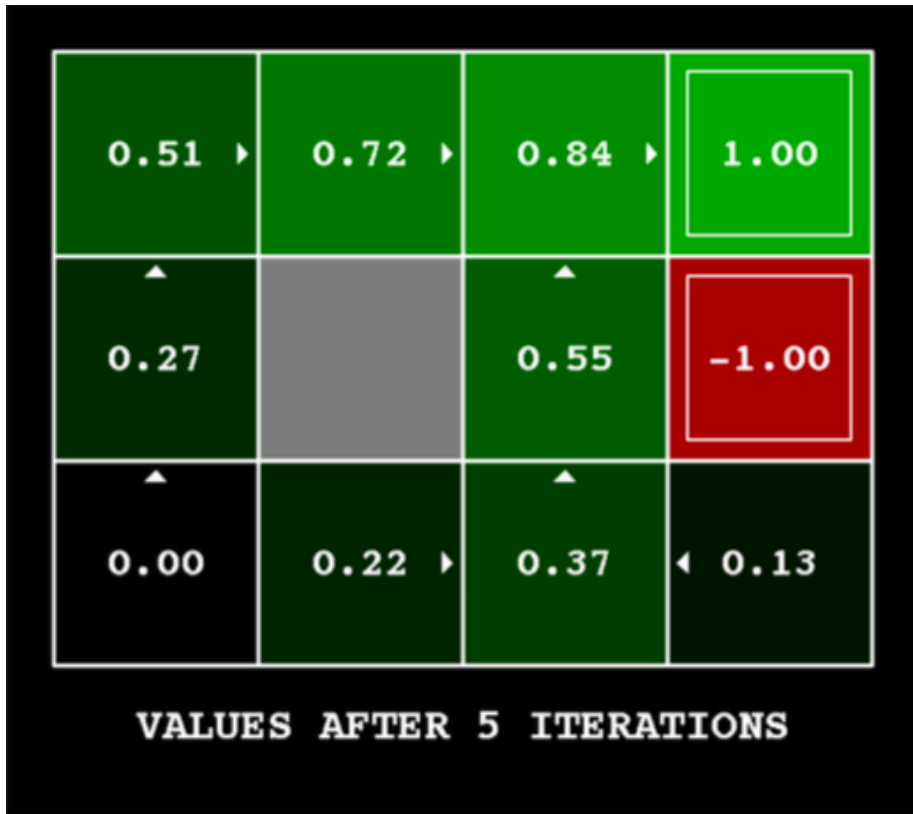VALUES AFTER 8 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s, a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 8 ITERATIONS



VALUES AFTER 9 ITERATIONS

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 10 ITERATIONS


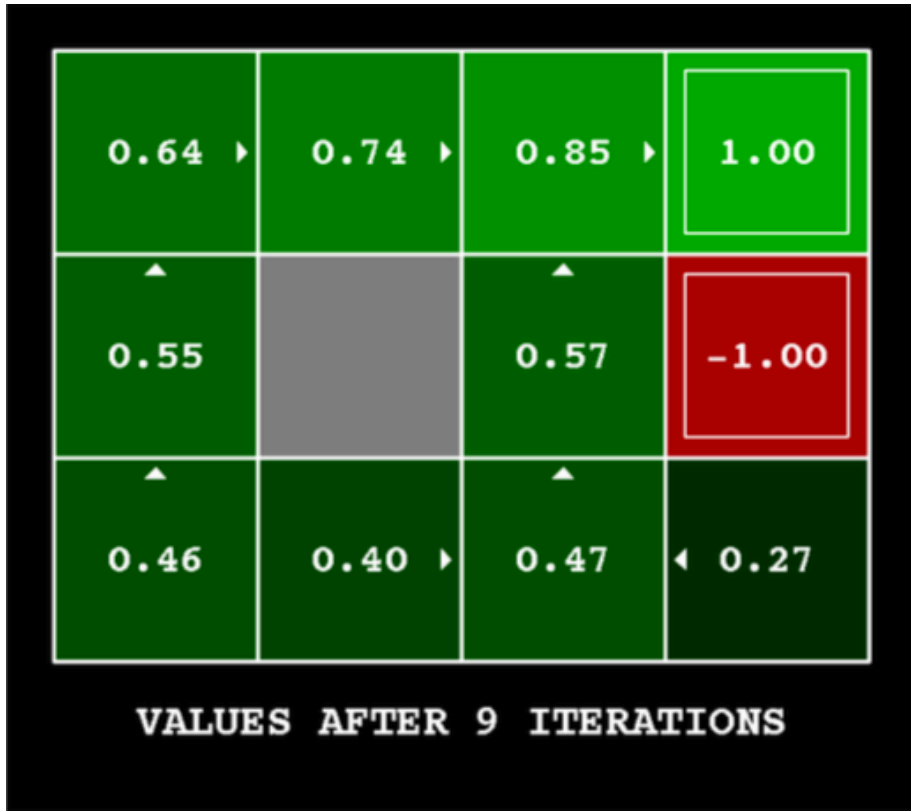
VALUES AFTER 11 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 11 ITERATIONS



VALUES AFTER 12 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# VI, GridWorld

$$V'(s) \leftarrow \max_{a \in A} \left[ r(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right], \quad \forall s$$



VALUES AFTER 12 ITERATIONS

VALUES AFTER 100 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9

# Policy Iteration

Don't stop in goal states in this Grid World

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | ■ | 0 | -100 |
| 0 | 0 | 0 | 0 |

Rewards

(Z. Kolter)

Initialize "up" everywhere

| | | | |
|---|---|---|---|
| 0.418 | 0.884 | 2.331 | 6.367 |
| 0.367 | ■ | -8.610 | -105.7 |
| -0.168 | -4.641 | -14.27 | -85.05 |

After 1 improvement step

pi0-> V0 -> pi1 -> V1

# Policy Iteration

Don't stop in goal states in this Grid World

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | ■ | 0 | -100 |
| 0 | 0 | 0 | 0 |

Rewards

(Z. Kolter)

| | | | |
|---|---|---|---|
| 5.414 | 6.248 | 7.116 | 8.634 |
| 4.753 | ■ | 2.881 | -102.7 |
| 2.251 | 1.977 | 1.849 | -8.701 |

After 2 improvement steps

pi0-> V0 -> pi1 -> V1 -> pi2 -> V2

# Policy Iteration

Don't stop in goal states in this Grid World

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 1 |
| 0 | ■ | 0 | -100 |
| 0 | 0 | 0 | 0 |

Rewards

| | | | |
|---|---|---|---|
| 5.470 | 6.313 | 7.190 | 8.669 |
| 4.803 | ■ | 3.347 | -96.67 |
| 4.161 | 3.654 | 3.222 | 1.526 |

After 3 improvement steps (converged!)

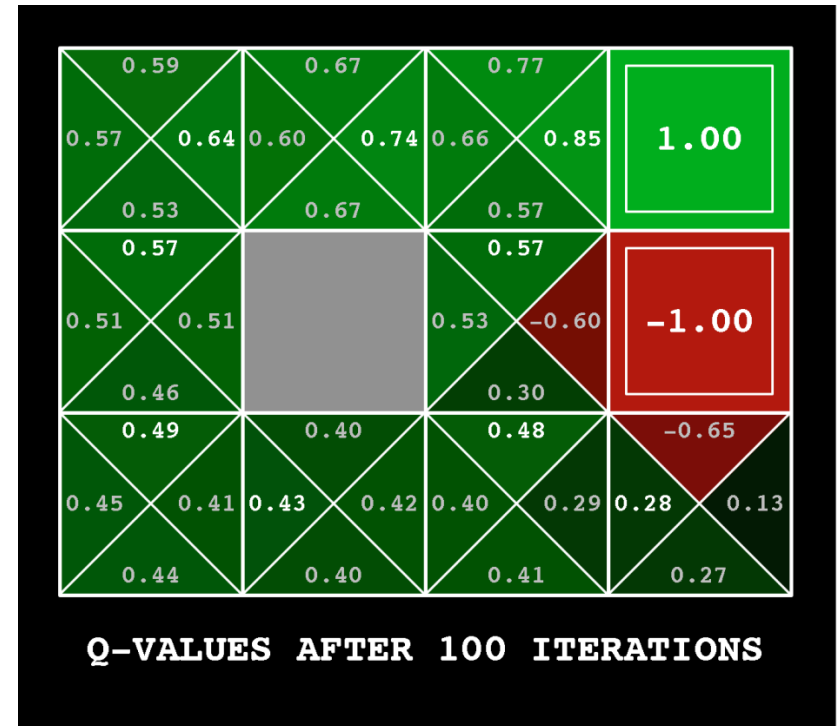pi0-> V0 -> pi1 -> V1 -> pi2 -> V2 -> pi3 -> V3

# GridWorld



(D. Klein, P. Abbeel)

- r(s,a) = 0, except states (1,4), (2,4). In these states get +1 or -1 when take ANY action. Then no more actions

- Bounce off obstacles. Actuator has 20% noise; e.g., w/ prob 0.1 goes L, prob 0.1 goes R when moving U

- Discounting 0.9 (r + 0.9 r + $0.9^2$ r + ...)

# Can also look at Q-Values



VALUES AFTER 100 ITERATIONS



Q-VALUES AFTER 100 ITERATIONS

(D. Klein, P. Abbeel)

Noise 0.2, Discount 0.9