

Sequential decisions under uncertainty

Markov Decision Processes (MDP)

Tomáš Svoboda

Vision for Robots and Autonomous Systems, Center for Machine Perception
Department of Cybernetics
Faculty of Electrical Engineering, Czech Technical University in Prague

April 1, 2020

slyším velmi dobře	<input type="checkbox"/>
slyším nic moc, ale dostatečně	<input type="checkbox"/> 1
nedostatečně	<input type="checkbox"/> 0

1 / 28

Sequential decisions under uncertainty

Markov Decision Processes (MDP)

Tomáš Svoboda

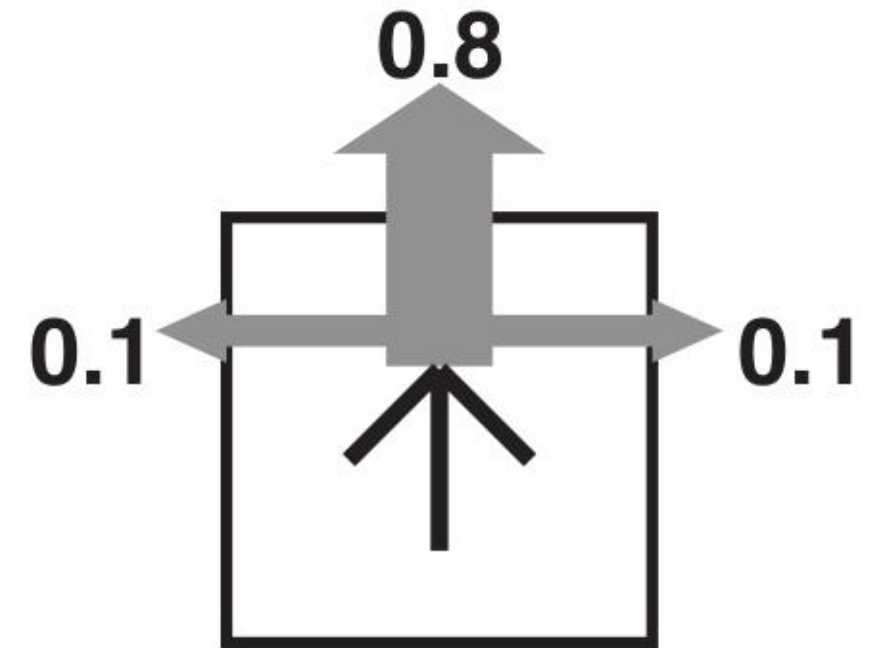
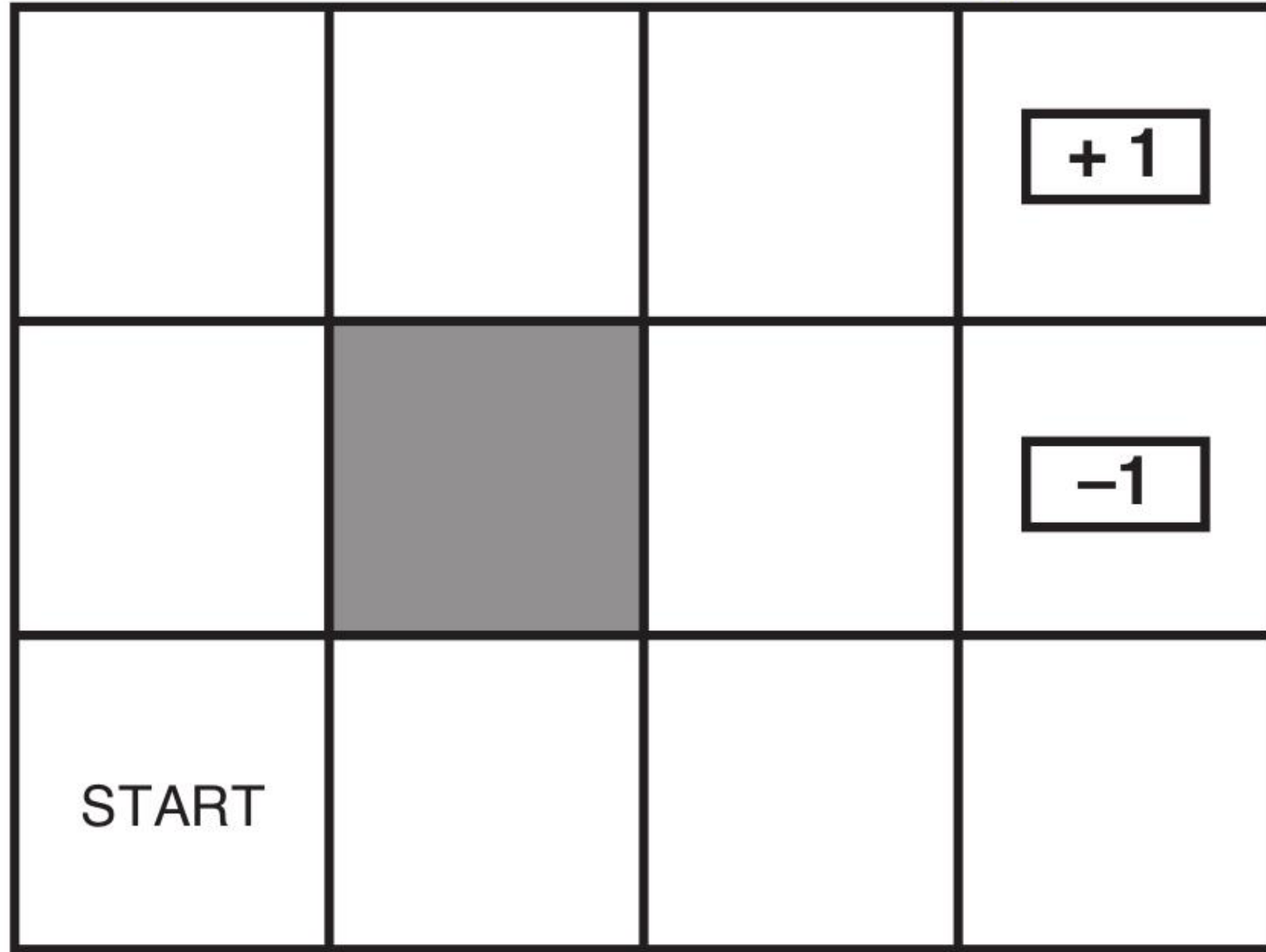
Vision for Robots and Autonomous Systems, Center for Machine Perception
Department of Cybernetics
Faculty of Electrical Engineering, Czech Technical University in Prague

April 1, 2020

slyším velmi dobře
řnám na serveru
i video pro off-line



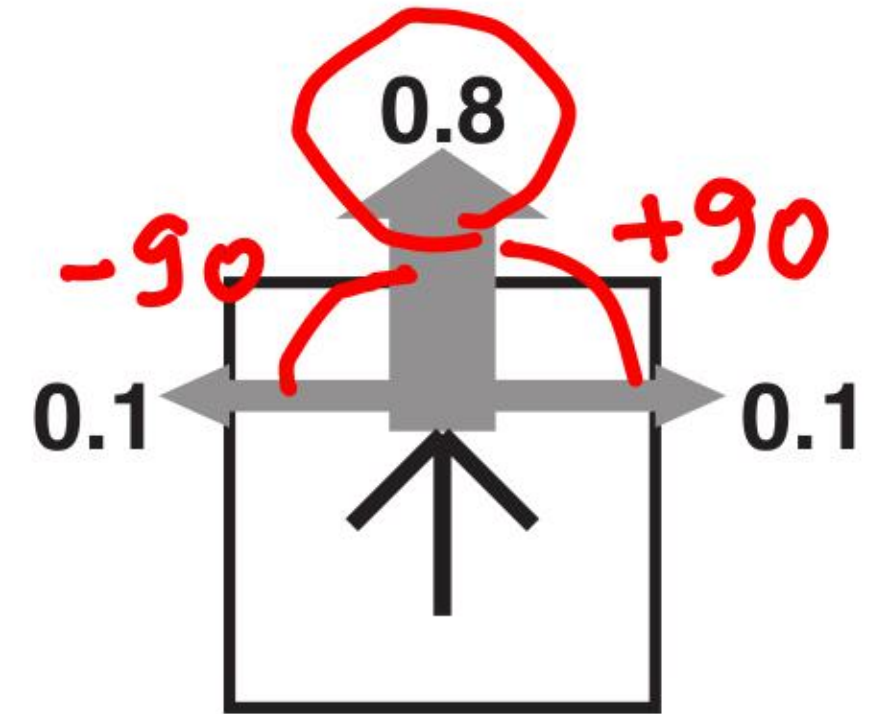
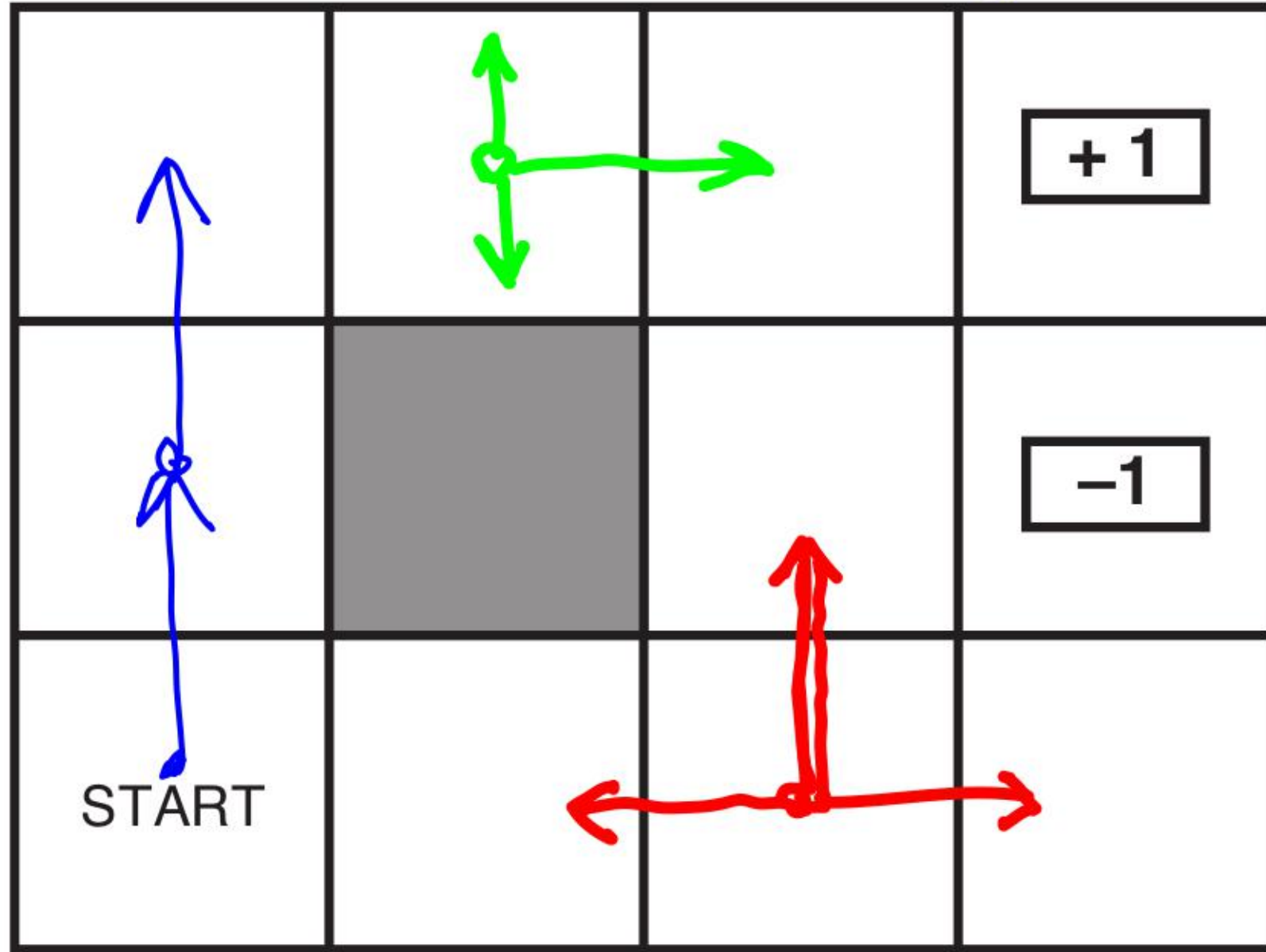
Unreliable actions in observable grid world



States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$

(Transition) Model $T(s, a, s') \equiv p(s'|s, a)$ = probability that a in s leads to s'

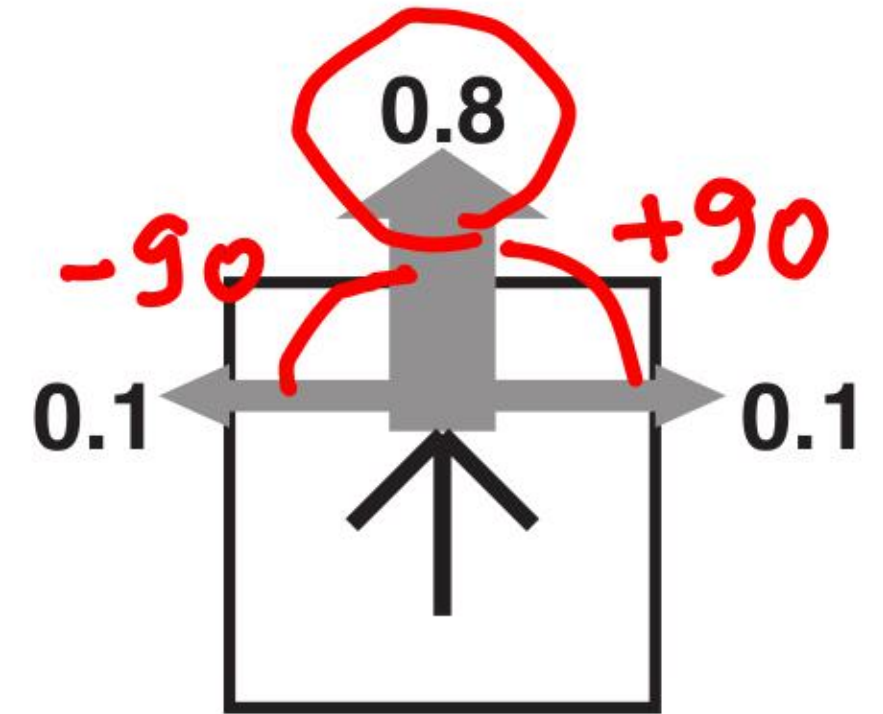
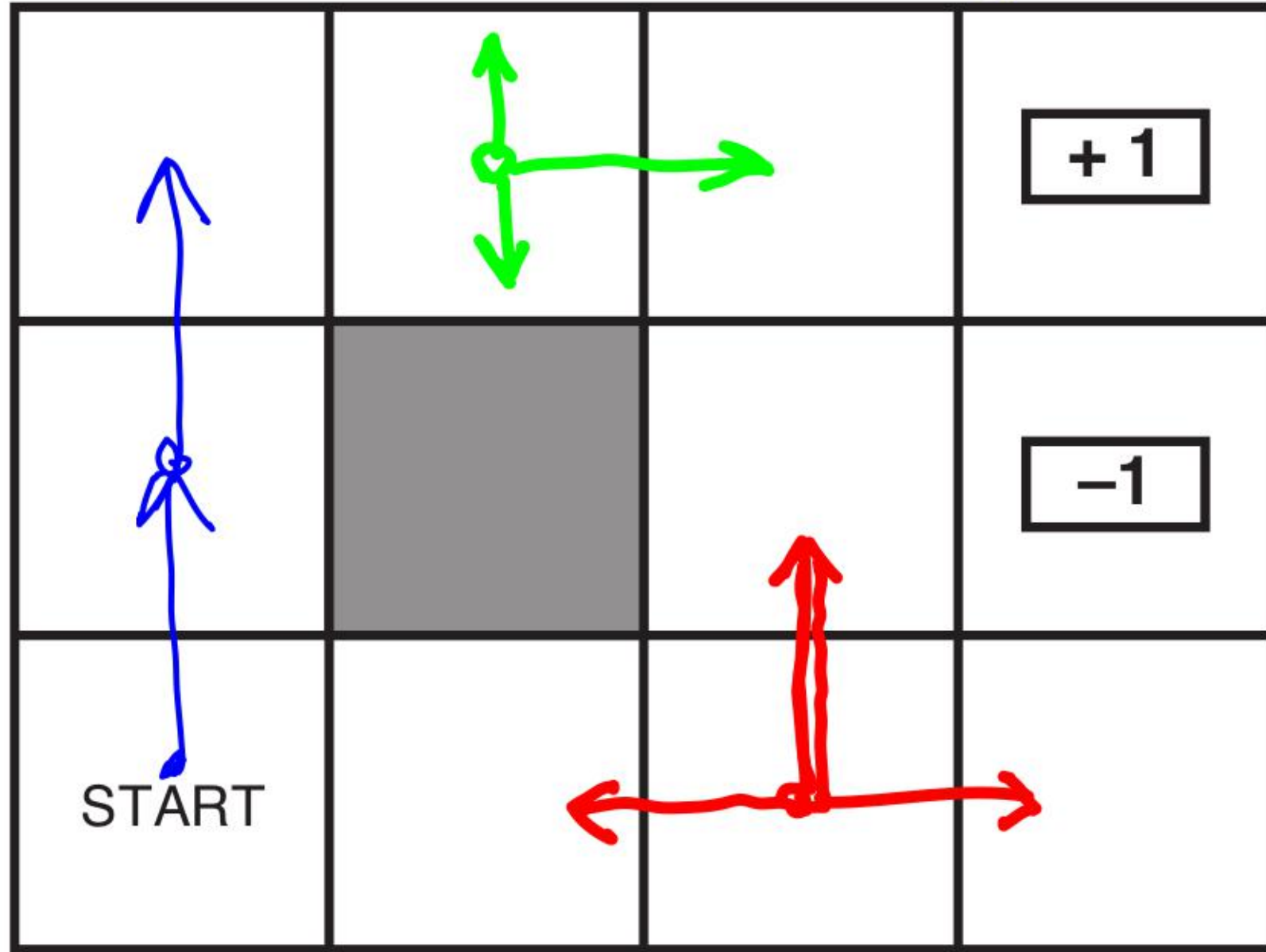
Unreliable actions in observable grid world



States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$

(Transition) Model $T(s, a, s') \equiv p(s'|s, a)$ = probability that a in s leads to s'

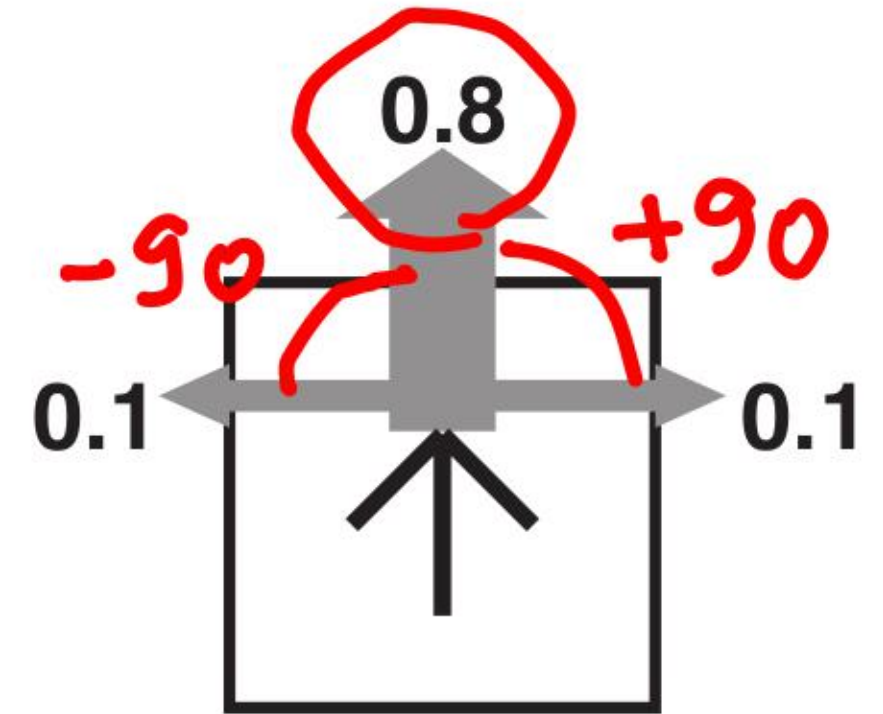
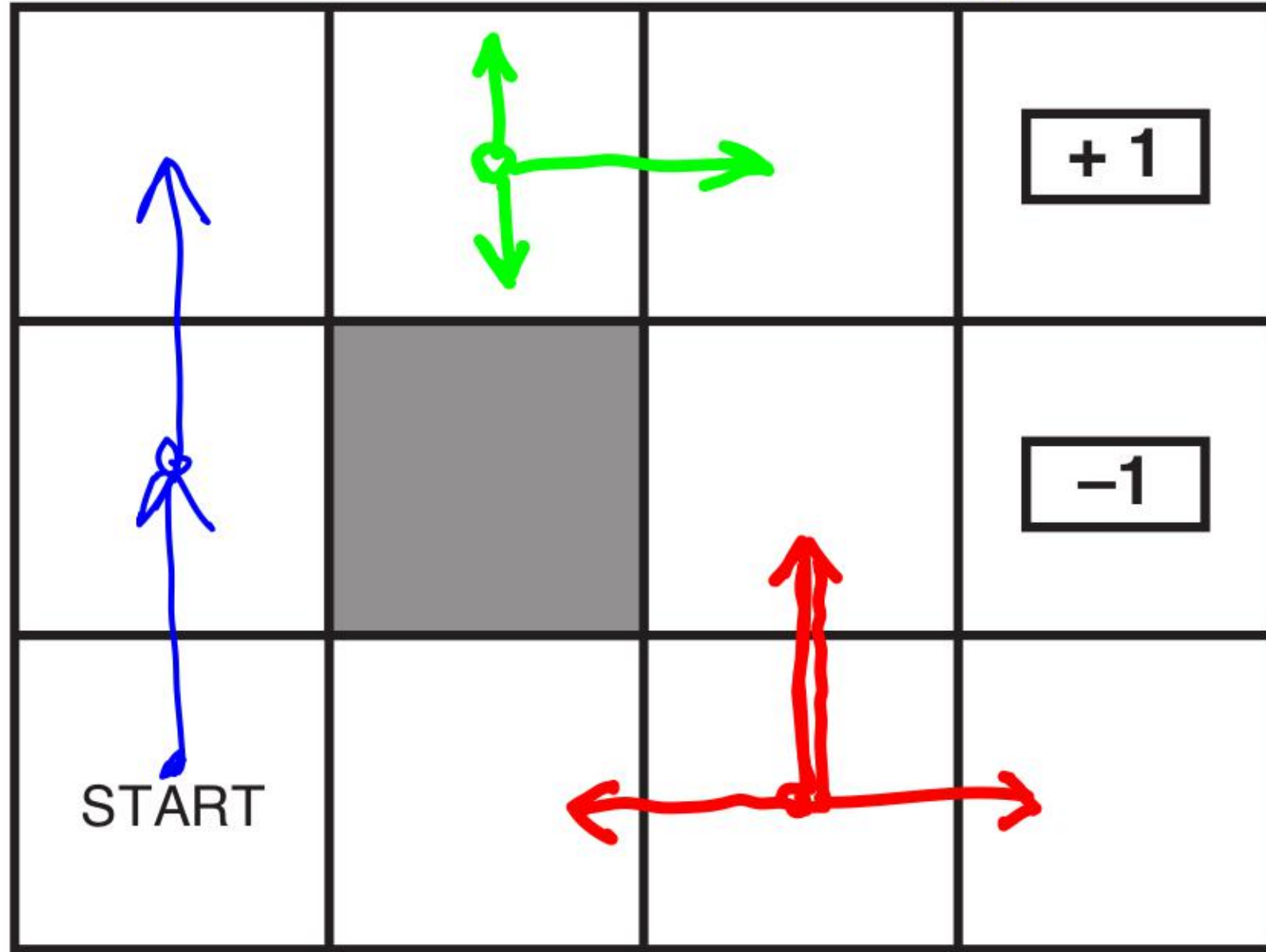
Unreliable actions in observable grid world



States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$

(Transition) Model $T(s, a, s') \equiv p(s'|s, a)$ = probability that a in s leads to s'

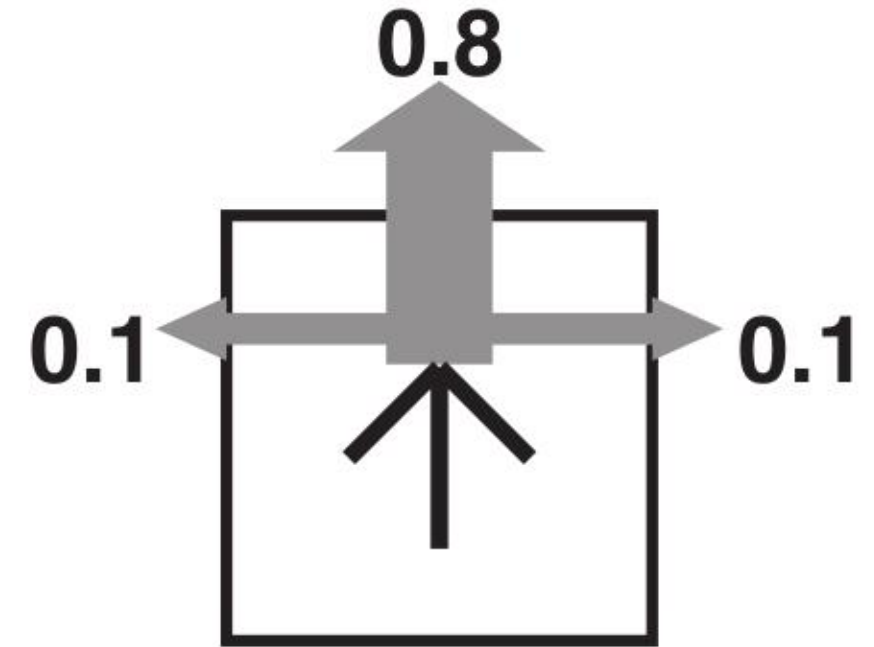
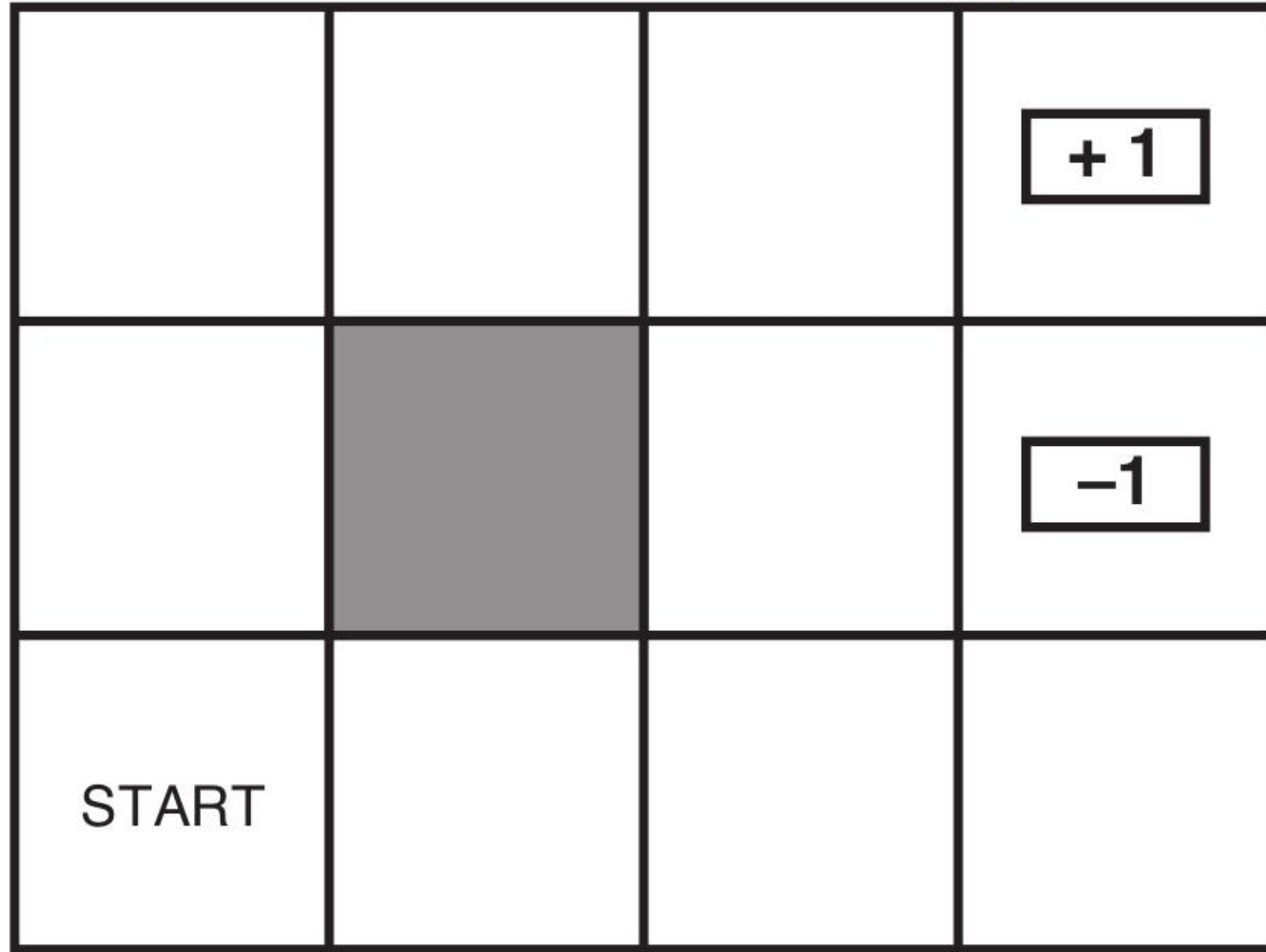
Unreliable actions in observable grid world



States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$

(Transition) Model $T(s, a, s') \equiv p(s'|s, a)$ = probability that a in s leads to s'

Unreliable actions in observable grid world



States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$

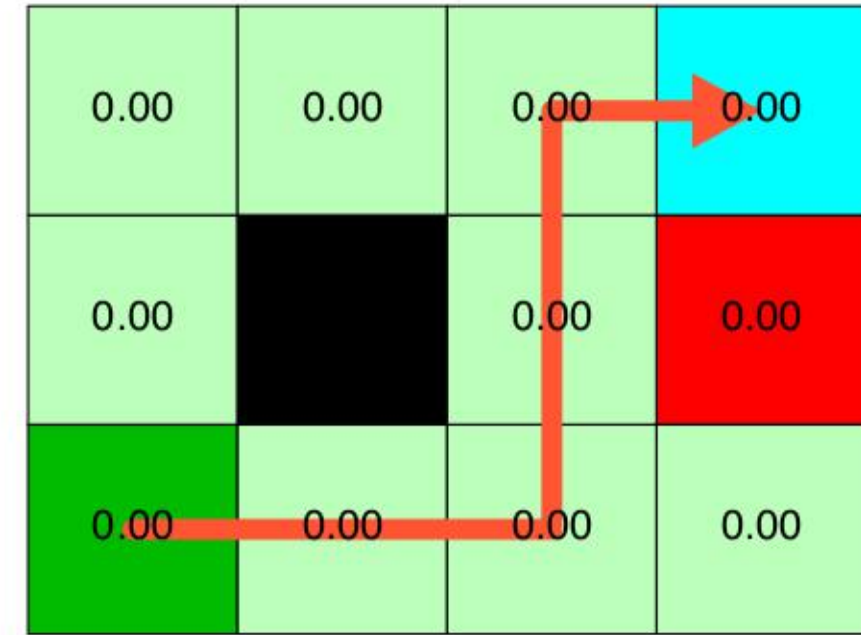
(Transition) Model $T(s, a, s') \equiv p(s'|s, a) =$ probability that a in s leads to s'

Unreliable actions



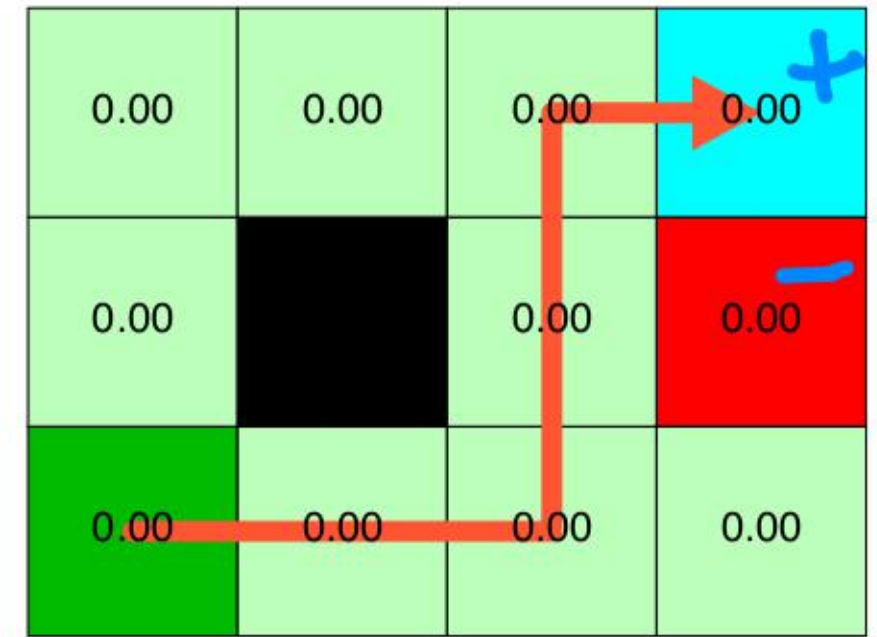
Plan? Policy

- ▶ In deterministic world: **Plan** – sequence of actions from **Start** to **Goal**.
- ▶ MDPs, we need a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- ▶ An action for each possible state. Why *each*?
- ▶ What is the *best* policy?



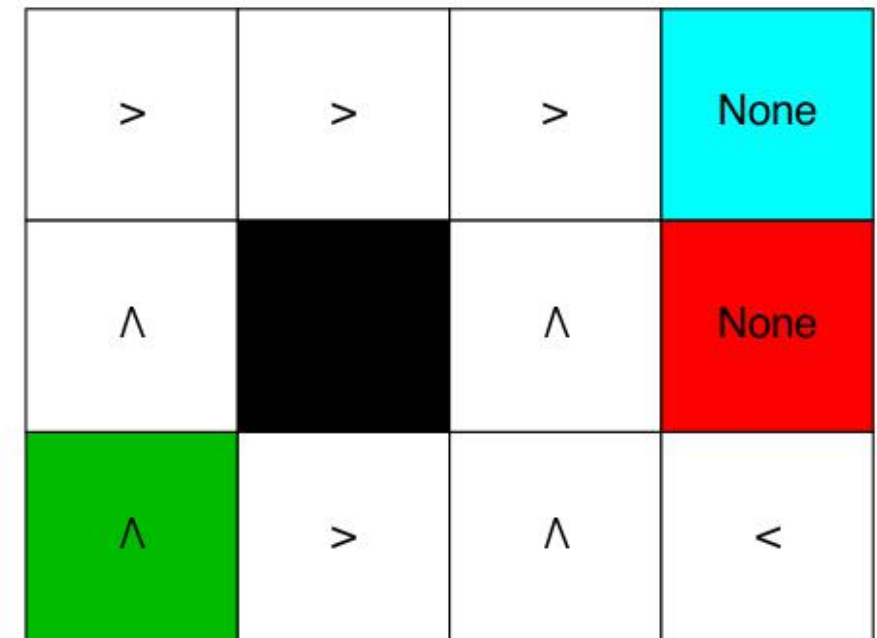
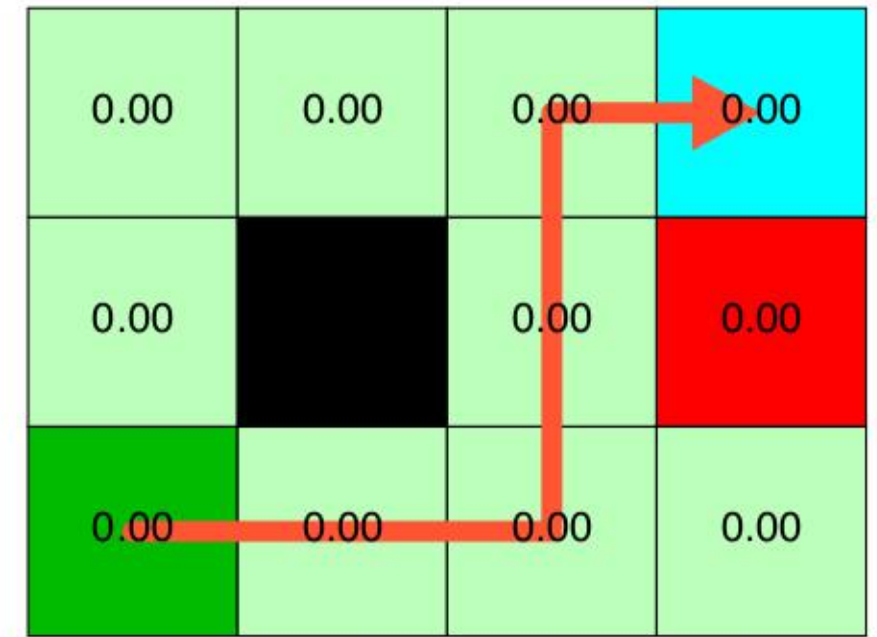
Plan? Policy

- ▶ In deterministic world: Plan – sequence of actions from Start to Goal.
- ▶ MDPs, we need a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- ▶ An action for each possible state. Why *each*?
- ▶ What is the *best* policy?



Plan? Policy

- ▶ In deterministic world: **Plan** – sequence of actions from **Start** to **Goal**.
- ▶ MDPs, we need a **policy** $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- ▶ An action for each possible state. *Why each?*
- ▶ What is the *best* policy?



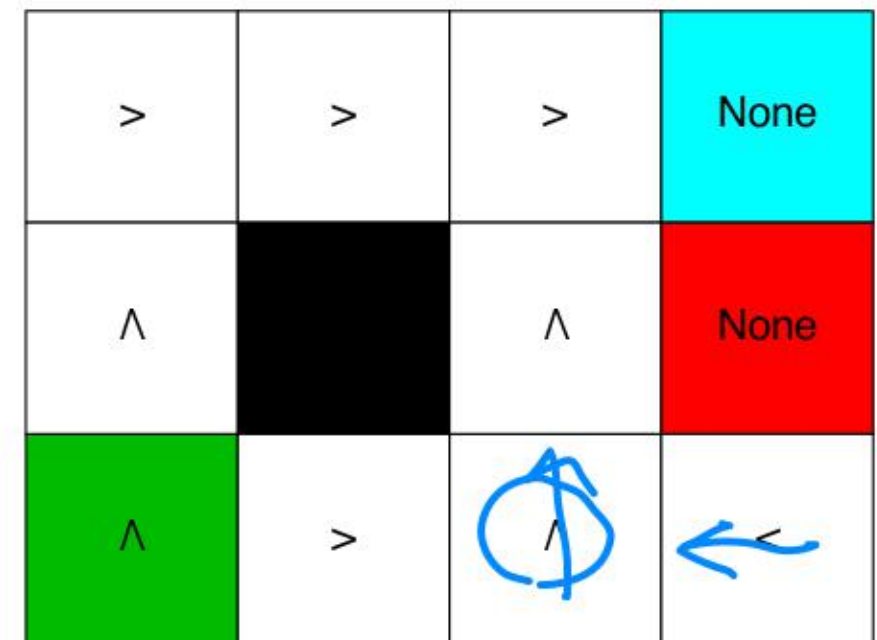
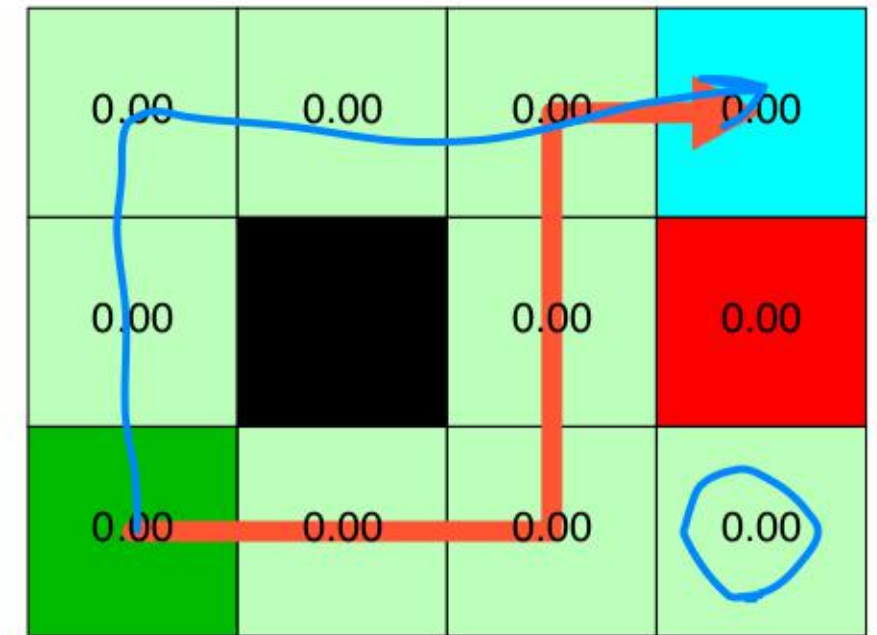
Plan? Policy

- ▶ In deterministic world: Plan – sequence of actions from Start to Goal.

SDP MDPs, we need a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$. *Strategie*

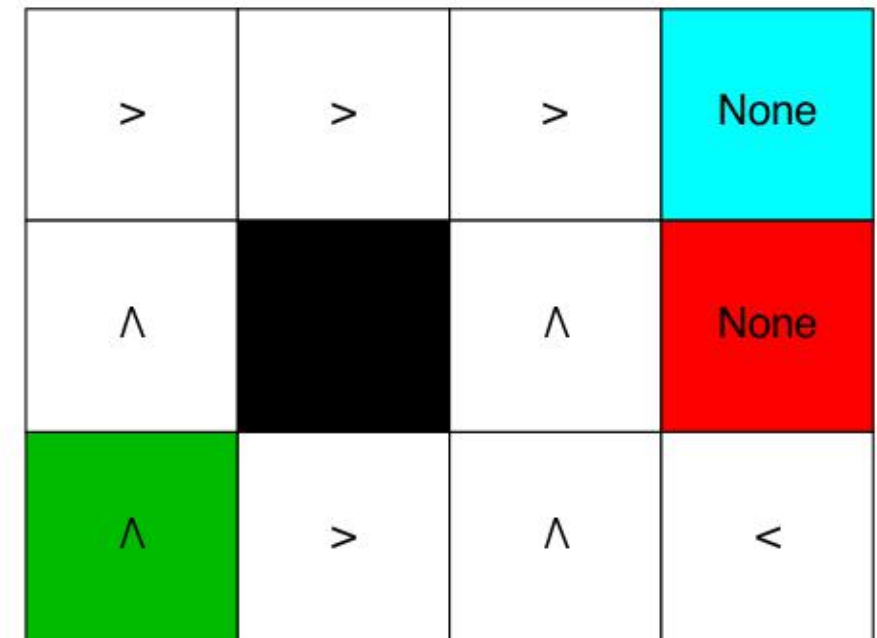
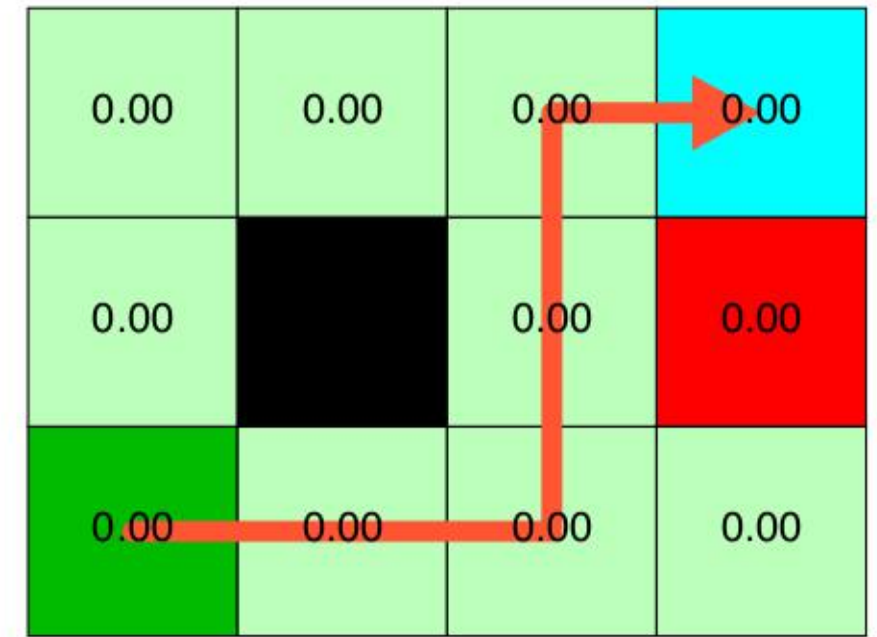
- ▶ An action for each possible state. Why each?

▶ What is the *best* policy?



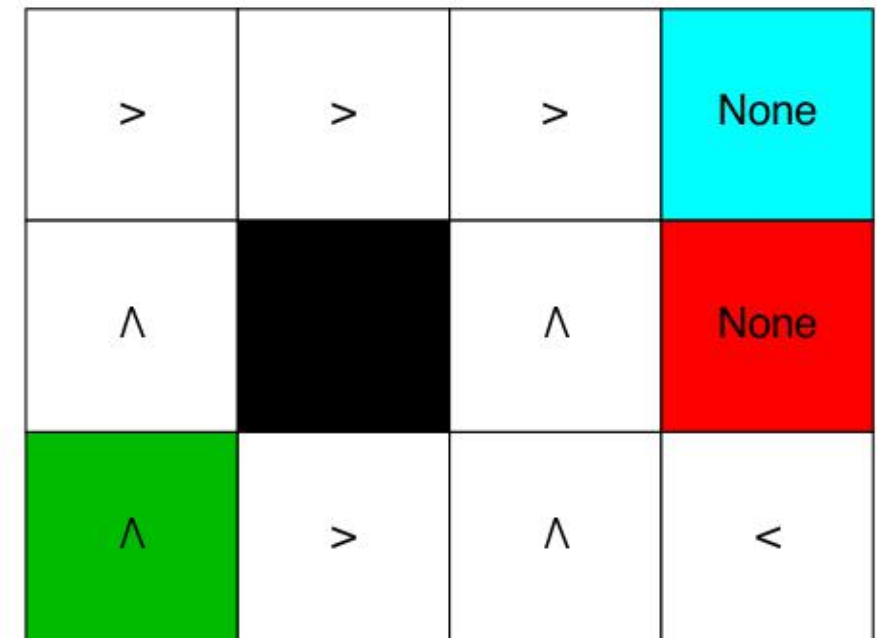
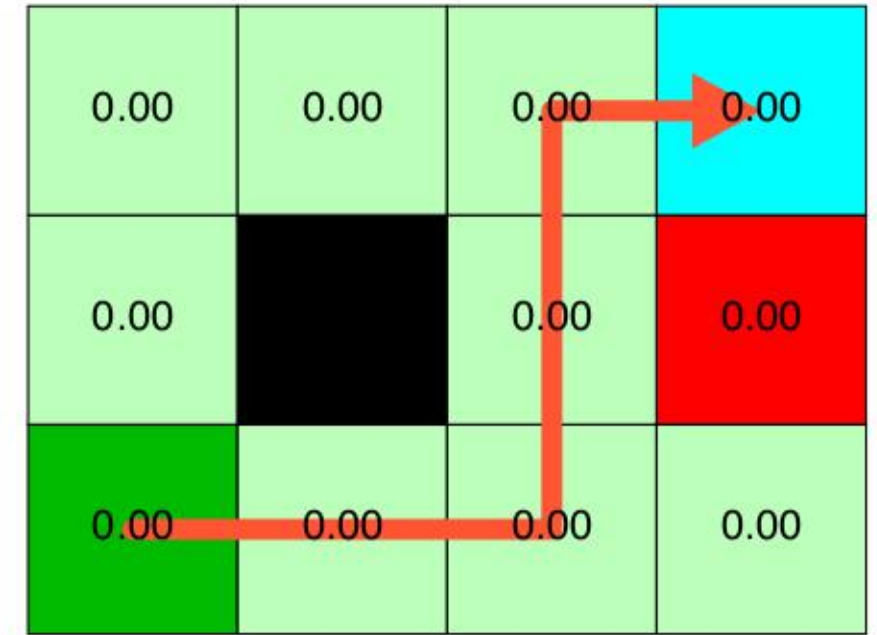
Plan? Policy

- ▶ In deterministic world: **Plan** – sequence of actions from **Start** to **Goal**.
- ▶ MDPs, we need a **policy** $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- ▶ An action for each possible state. *Why each?*
- ▶ What is the *best* policy?



Plan? Policy

- ▶ In deterministic world: Plan – sequence of actions from Start to Goal.
- ▶ MDPs, we need a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- ▶ An action for each possible state. Why *each*?
- ▶ What is the best policy?



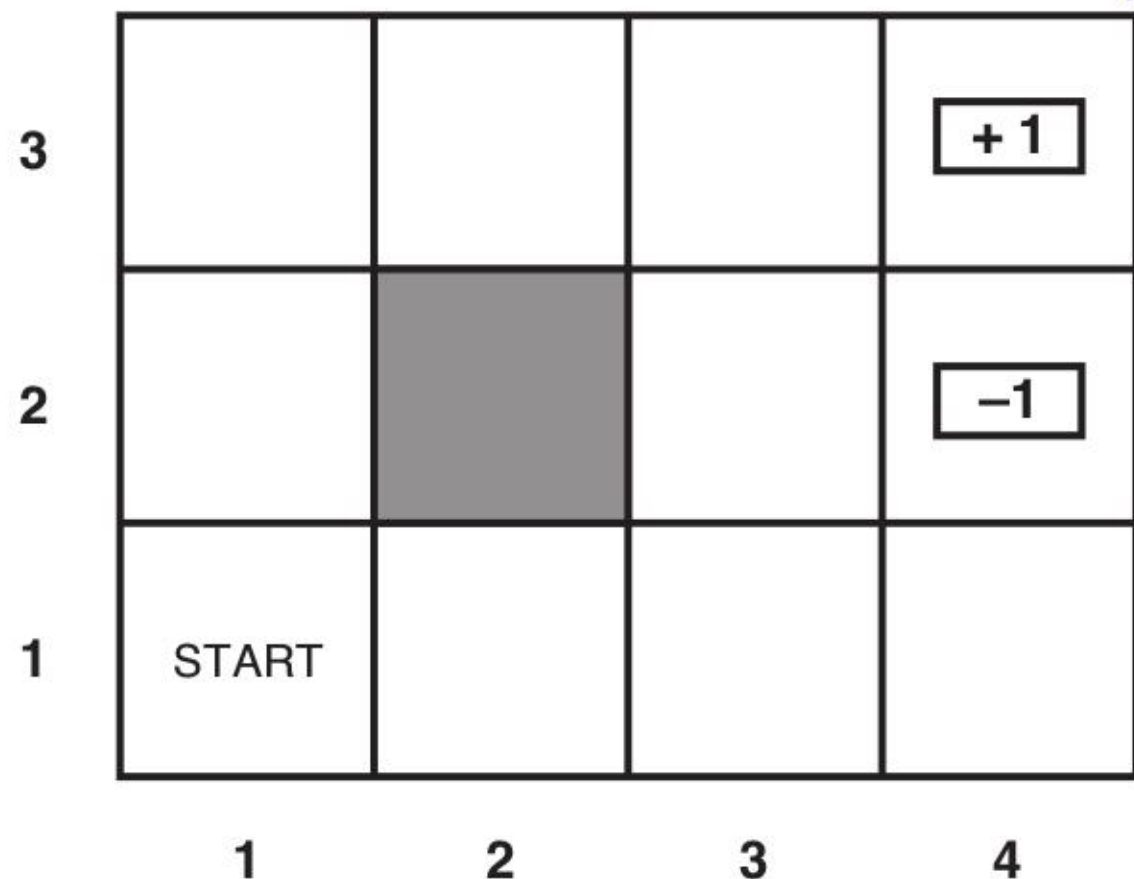
Rewards

-0.04	-0.04	-0.04	1.00
-0.04		-0.04	-1.00
-0.04	-0.04	-0.04	-0.04

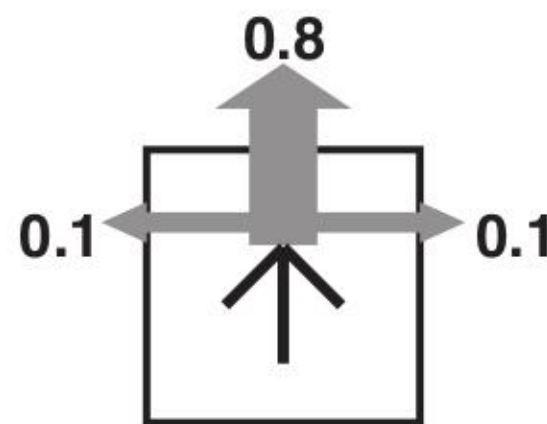
Reward : Robot/Agent takes an action a and it is **immediately** rewarded.

Reward function $r(s)$ (or $r(s, a)$, $r(s, a, s')$)
= $\begin{cases} -0.04 & \text{(small penalty) for nonterminal states} \\ \pm 1 & \text{for terminal states} \end{cases}$

Markov Decision Processes (MDPs)



(a)



(b)

States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$

Model $T(s, a, s') \equiv p(s'|s, a)$ = probability that a in s leads to s'

Reward function $r(s)$ (or $r(s, a)$, $r(s, a, s')$)

$$= \begin{cases} -0.04 & \text{(small penalty) for nonterminal states} \\ \pm 1 & \text{for terminal states} \end{cases}$$

Markovian property

- ▶ Given the present state, the future and the past are independent.
- ▶ MDP: Markov means action depends only on the current state.
- ▶ In search: successor function (transition model) depends on the current state only.

Markovian property

- ▶ Given the present state, the future and the past are independent.
- ▶ MDP: Markov means action depends only on the current state.
- ▶ In search: successor function (transition model) depends on the current state only.

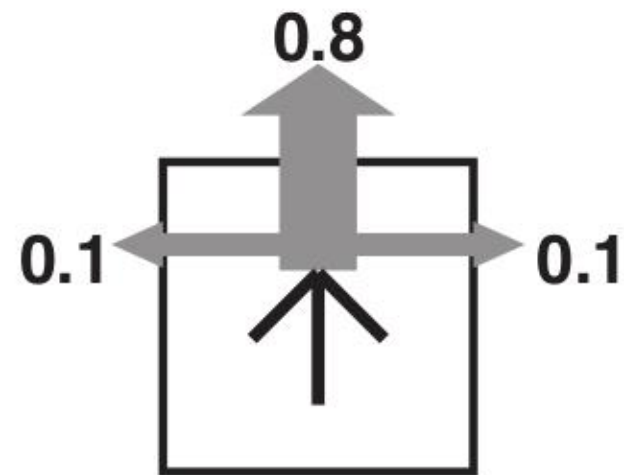
<table border="1"> <tbody> <tr> <td>></td> <td>></td> <td>></td> <td>1.00</td> </tr> <tr> <td>^</td> <td></td> <td>^</td> <td>-1.00</td> </tr> <tr> <td>^</td> <td><</td> <td><</td> <td><</td> </tr> </tbody> </table> <p style="text-align: center;">A</p> <p style="text-align: center;">$r(s) = \{-2, 1, -1\}$ a</p>	>	>	>	1.00	^		^	-1.00	^	<	<	<	<table border="1"> <tbody> <tr> <td>></td> <td>></td> <td>></td> <td>1.00</td> </tr> <tr> <td>^</td> <td></td> <td><</td> <td>-1.00</td> </tr> <tr> <td>^</td> <td><</td> <td><</td> <td>v</td> </tr> </tbody> </table> <p style="text-align: center;">B</p> <p style="text-align: center;">$r(s) = \{-0.04, 1, -1\}$ b</p>	>	>	>	1.00	^		<	-1.00	^	<	<	v	<table border="1"> <tbody> <tr> <td>></td> <td>></td> <td>></td> <td>1.00</td> </tr> <tr> <td>^</td> <td></td> <td>></td> <td>-1.00</td> </tr> <tr> <td>></td> <td>></td> <td>></td> <td>^</td> </tr> </tbody> </table> <p style="text-align: center;">C</p> <p style="text-align: center;">$r(s) = \{-0.01, 1, -1\}$ c</p>	>	>	>	1.00	^		>	-1.00	>	>	>	^
>	>	>	1.00																																			
^		^	-1.00																																			
^	<	<	<																																			
>	>	>	1.00																																			
^		<	-1.00																																			
^	<	<	v																																			
>	>	>	1.00																																			
^		>	-1.00																																			
>	>	>	^																																			

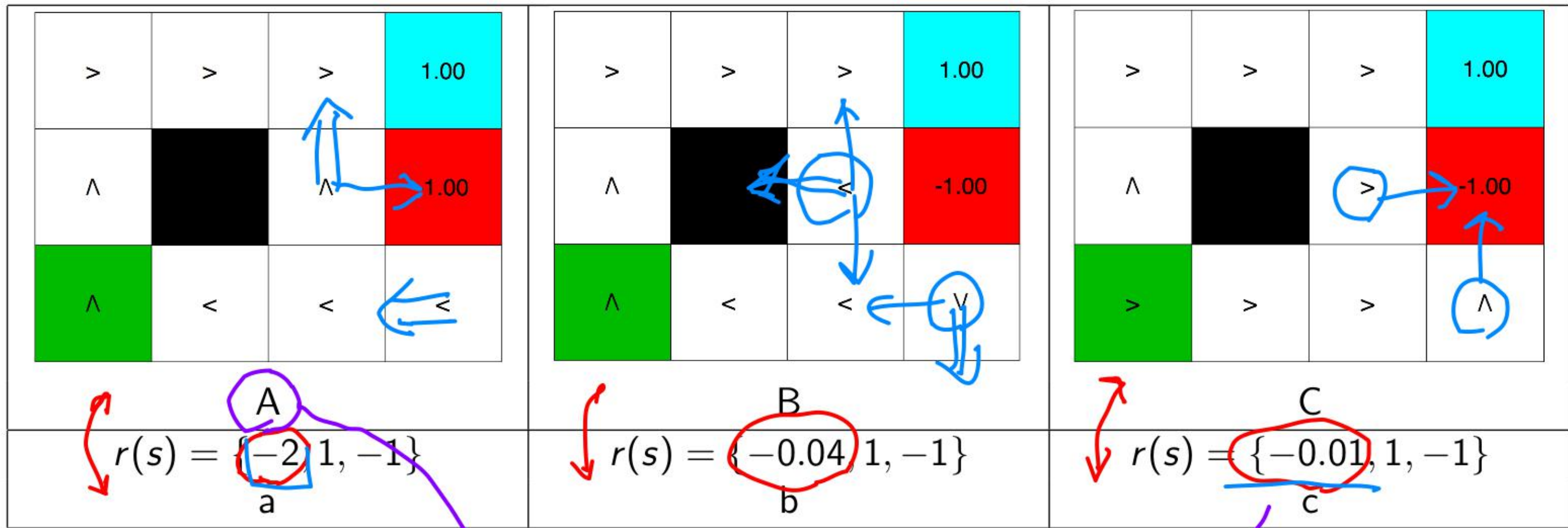
A: A-a, B-b, C-c

B: A-b, B-a, C-c

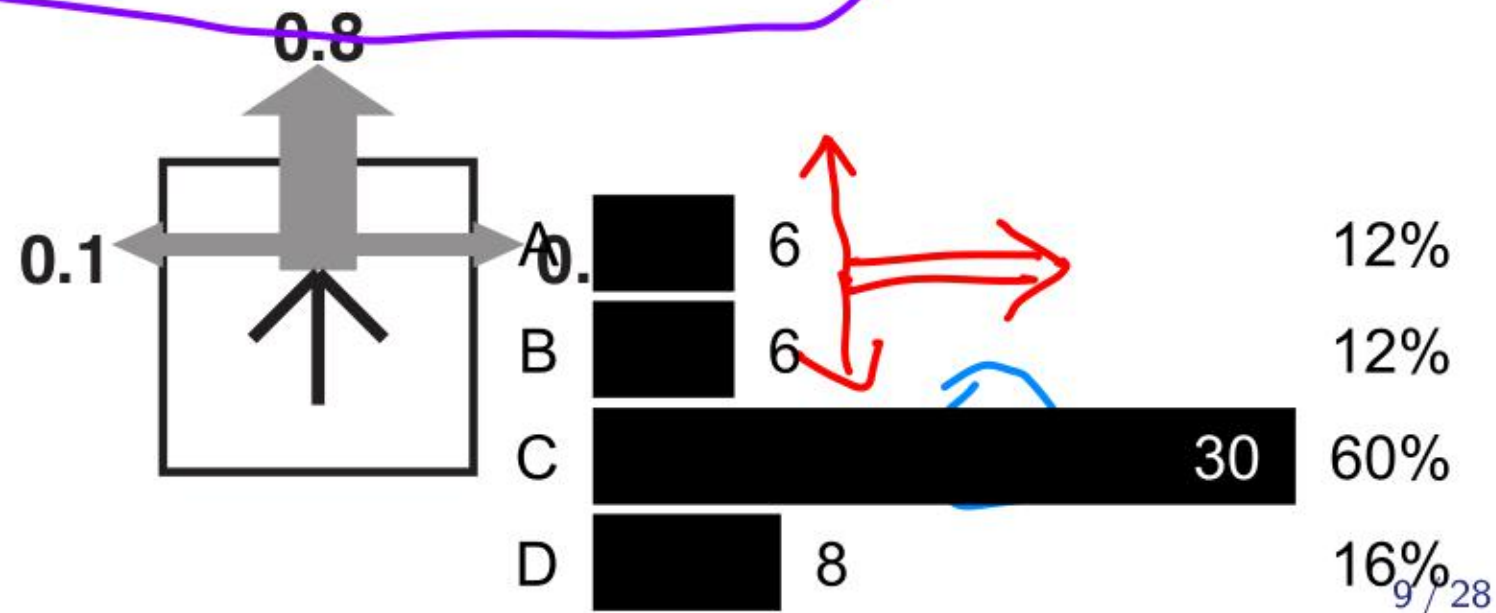
C: A-b, B-c, C-a

D: A-c, B-a, C-b





- A: A-a, B-b, C-c
- B: A-b, B-a, C-c
- C: A-b, B-c, C-a
- D: A-c, B-a, C-b



Utilities of sequences

- ▶ State reward at time/step t , R_t .
- ▶ State at time t , S_t . State sequence $[S_0, S_1, S_2, \dots,]$

Typically, consider stationary preferences on reward sequences:

$$[R, R_1, R_2, R_3, \dots] \succ [R, R'_1, R'_2, R'_3, \dots] \Leftrightarrow [R_1, R_2, R_3, \dots] \succ [R'_1, R'_2, R'_3, \dots]$$

If stationary preferences :

Utility (h -history)

$$U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \dots$$

If the horizon is finite - limited number of steps - preferences are nonstationary (depends on how many steps left).

Utilities of sequences



- ▶ State reward at time/step t , R_t .
- ▶ State at time t , S_t . State sequence $[S_0, S_1, S_2, \dots,]$

Typically, consider stationary preferences on reward sequences:

$$[R, R_1, R_2, R_3, \dots] \succ [R, R'_1, R'_2, R'_3, \dots] \Leftrightarrow [R_1, R_2, R_3, \dots] \succ [R'_1, R'_2, R'_3, \dots]$$

If stationary preferences :

Utility (h -history)

$$U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \dots$$

If the horizon is finite - limited number of steps - preferences are nonstationary (depends on how many steps left).

Utilities of sequences

- ▶ State reward at time/step t , R_t .
- ▶ State at time t , S_t . State sequence $[S_0, S_1, S_2, \dots,]$

Typically, consider **stationary preferences** on reward sequences:

$$[R, R_1, R_2, R_3, \dots] \succ [R, R'_1, R'_2, R'_3, \dots] \Leftrightarrow [R_1, R_2, R_3, \dots] \succ [R'_1, R'_2, R'_3, \dots]$$

If stationary preferences :

Utility (h -history)

$$U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \dots$$

If the horizon is finite - limited number of steps - preferences are nonstationary (depends on how many steps left).

Utilities of sequences

- ▶ State reward at time/step t , R_t .
- ▶ State at time t , S_t . State sequence $[S_0, S_1, S_2, \dots,]$

Typically, consider **stationary preferences** on reward sequences:

$$[R, R_1, R_2, R_3, \dots] \succ [R, R'_1, R'_2, R'_3, \dots] \Leftrightarrow [R_1, R_2, R_3, \dots] \succ [R'_1, R'_2, R'_3, \dots]$$

If **stationary preferences** :

Utility (h -history)

$$U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \dots$$

If the horizon is finite - limited number of steps - preferences are **nonstationary** (depends on how many steps left).

Utilities of sequences

- ▶ State reward at time/step t , R_t .
- ▶ State at time t , S_t . State sequence $[S_0, S_1, S_2, \dots,]$

Typically, consider **stationary preferences** on reward sequences:

$$[R, R_1, R_2, R_3, \dots] \succ [R, R'_1, R'_2, R'_3, \dots] \Leftrightarrow [R_1, R_2, R_3, \dots] \succ [R'_1, R'_2, R'_3, \dots]$$

If **stationary preferences** :

Utility (h -history)

$$U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \dots$$

If the horizon is finite - limited number of steps - preferences are **nonstationary** (depends on how many steps left).

Returns and Episodes

- ▶ Executing policy - sequence of states and **rewards**.
- ▶ **Episode** starts at t , ends at T (ending in a terminal state).
- ▶ **Return** (Utility) of the episode (policy execution)

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

Returns and Episodes

- ▶ Executing policy → sequence of states and **rewards**.
- ▶ Episode starts at t , ends at T (ending in a terminal state).
- ▶ Return (Utility) of the episode (policy execution)

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

Comparing policies; Finite vs infinite horizon

Problem: Infinite lifetime \Rightarrow additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- ▶ Discounted return , $\gamma < 1, R_t \leq R_{\max}$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \frac{R_{\max}}{1-\gamma}$$

- ▶ Absorbing (terminal) state.

Returns are successive steps related to each other

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Comparing policies; Finite vs infinite horizon

Problem: Infinite lifetime \Rightarrow additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- ▶ Discounted return, $\gamma < 1$, $R_t \leq R_{\max}$

$$\underline{G}_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 \underline{R}_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \frac{R_{\max}}{1 - \gamma}$$

- ▶ Absorbing (terminal) state.

Returns are successive steps related to each other

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Comparing policies; Finite vs infinite horizon

Problem: Infinite lifetime \Rightarrow additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- ▶ Discounted return , $\gamma < 1, R_t \leq R_{\max}$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \frac{R_{\max}}{1-\gamma}$$

- ▶ Absorbing (terminal) state.

Returns are successive steps related to each other

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma^1 R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Comparing policies; Finite vs infinite horizon

Problem: Infinite lifetime \Rightarrow additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- ▶ Discounted return , $\gamma < 1, R_t \leq R_{\max}$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \frac{R_{\max}}{1-\gamma}$$

- ▶ Absorbing (terminal) state.

Returns are successive steps related to each other

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma^1 R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= \underline{R_{t+1}} + \gamma G_{t+1} \end{aligned}$$

MDPs recap

Markov decision processes (MDPs):

- ▶ Set of states \mathcal{S}
- ▶ Set of actions \mathcal{A}
- ▶ Transitions $p(s'|s, a)$ or $T(s, a, s')$
- ▶ Reward function $r(s, a, s')$; and discount γ

MDP quantities:

- ▶ (deterministic) Policy $\pi(s)$ – choice of action for each state
- ▶ Return (Utility) of an episode (sequence) – sum of (discounted) rewards.

MDPs recap

Markov decision processes (MDPs):

- ▶ Set of states \mathcal{S}
- ▶ Set of actions \mathcal{A}
- ▶ Transitions $p(s'|s, a)$ or $T(s, a, s')$
- ▶ Reward function $r(s, a, s')$; and discount γ

MDP quantities:

- ▶ (deterministic) Policy $\pi(s)$ – choice of action for each state
- ▶ Return (Utility) of an episode (sequence) – sum of (discounted) rewards.

Value functions

- ▶ Executing policy $\pi \rightarrow$ sequence of states (and rewards).
- ▶ Utility of a state sequence.
- ▶ But actions are unreliable - environment is stochastic.
- ▶ **Expected return** of a policy π .

Starting at time t , i.e. S_t ,

$$U^\pi(S_t) = \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right]$$

Value function

$$v^\pi(s) = \mathbb{E}^\pi [G_t \mid S_t = s] = \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Action-value function (q-function)

$$q^\pi(s, a) = \mathbb{E}^\pi [G_t \mid S_t = s, A_t = a] = \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

Value functions

- ▶ Executing policy $\pi \rightarrow$ sequence of states (and rewards).
- ▶ Utility of a state sequence.
- ▶ But actions are unreliable - environment is stochastic.
- ▶ Expected return of a policy π .

Starting at time t , i.e. S_t ,

$$U^\pi(S_t) = \mathbb{E}^\pi \left[\underbrace{\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}} \right]$$

Value function

$$v^\pi(s) = \mathbb{E}^\pi [G_t \mid S_t = s] = \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Action-value function (q-function)



$s \in \mathcal{N}$ $s \in \mathcal{S}$

$$q^\pi(s, a) = \mathbb{E}^\pi [G_t \mid S_t = s, A_t = a] = \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

Optimal policy π^* , and optimal value $v^*(s)$

$v^*(s)$ = expected (discounted) sum of rewards (until termination) assuming *optimal* actions.

	0	1	2	3
0	0.88	0.92	0.96	1.00
1	0.84		0.92	-1.00
2	0.80	0.84	0.88	0.84
	0	1	2	3

	0	1	2	3
0	>	>	>	None
1	^		^	None
2	^	>	^	<
	0	1	2	3

$$r(s) = \{-0.04, 1, -1\}, \gamma = 0.999999, \epsilon = 0.03, \text{ Robot } \textit{deterministic}$$

Optimal policy π^* , and optimal value $v^*(s)$

$v^*(s)$ = expected (discounted) sum of rewards (until termination) assuming *optimal* actions.

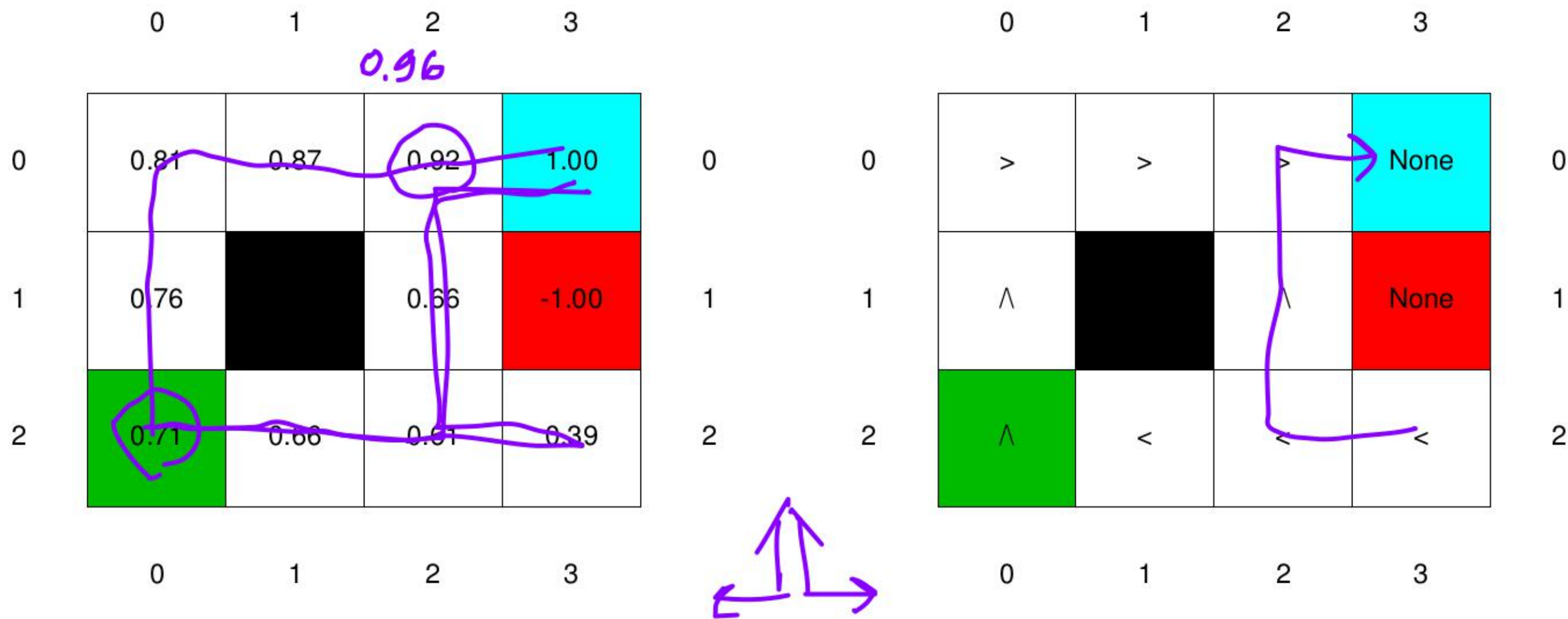
	0	1	2	3
0	0.81	0.87	0.92	1.00
1	0.76		0.66	-1.00
2	0.71	0.66	0.61	0.39
	0	1	2	3

	0	1	2	3
0	>	>	>	None
1	^		^	None
2	^	<	<	<
	0	1	2	3

$$r(s) = \{-0.04, 1, -1\}, \gamma = 0.9999999, \epsilon = 0.03, \text{ Robot } \textit{non-deterministic}$$

Optimal policy π^* , and optimal value $v^*(s)$

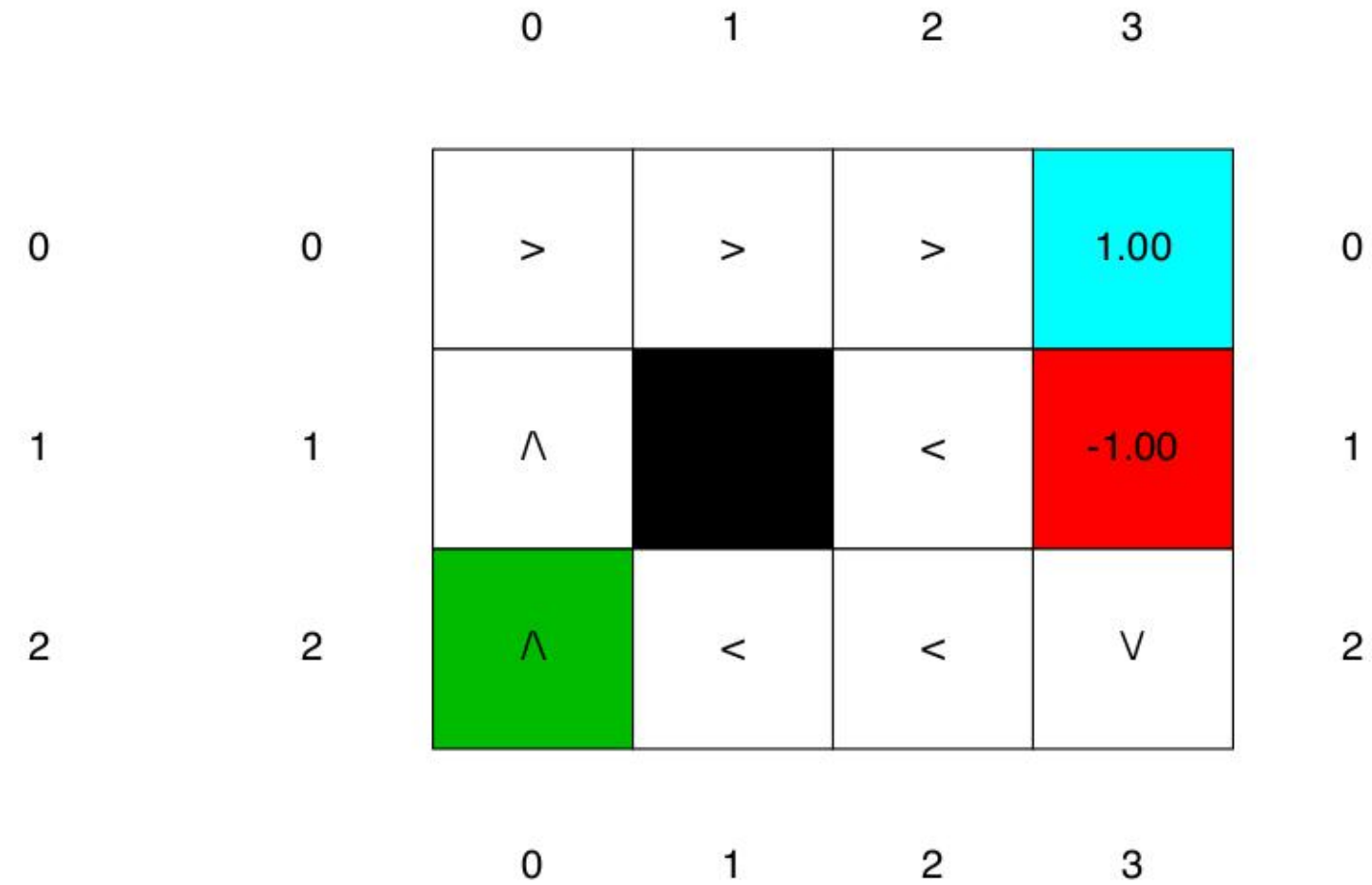
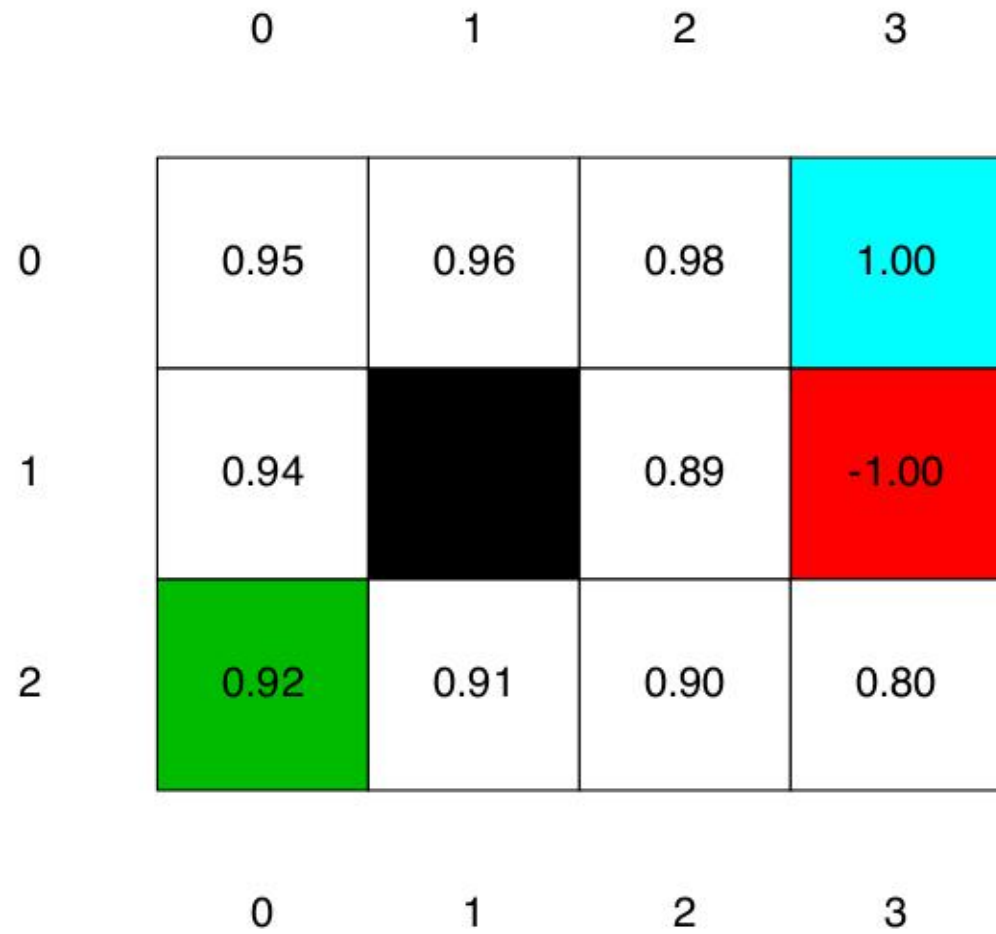
$v^*(s)$ = expected (discounted) sum of rewards (until termination) assuming *optimal* actions.



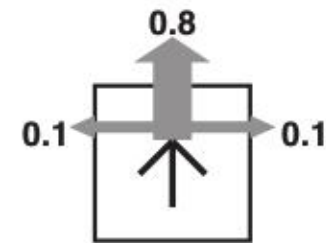
$r(s) = \{-0.04, 1, -1\}$, $\gamma = 0.9999999$, $\epsilon = 0.03$, Robot non-deterministic

Optimal policy π^* , and optimal value $v^*(s)$

$v^*(s)$ = expected (discounted) sum of rewards (until termination) assuming *optimal* actions.



$$r(s) = \{-0.01, 1, -1\}, \gamma = 0.9999999, \epsilon = 0.03,$$



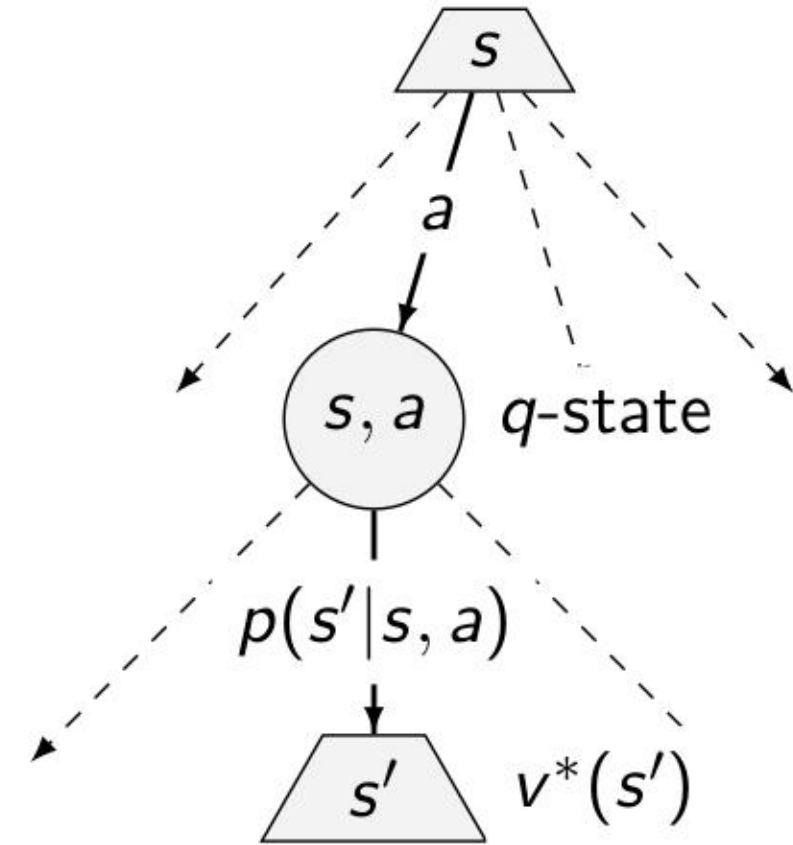
MDP search tree

The value of a q -state (s, a) :

$$q^*(s, a) = \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v^*(s')]$$

The value of a state s :

$$v^*(s) = \max_a q^*(s, a)$$



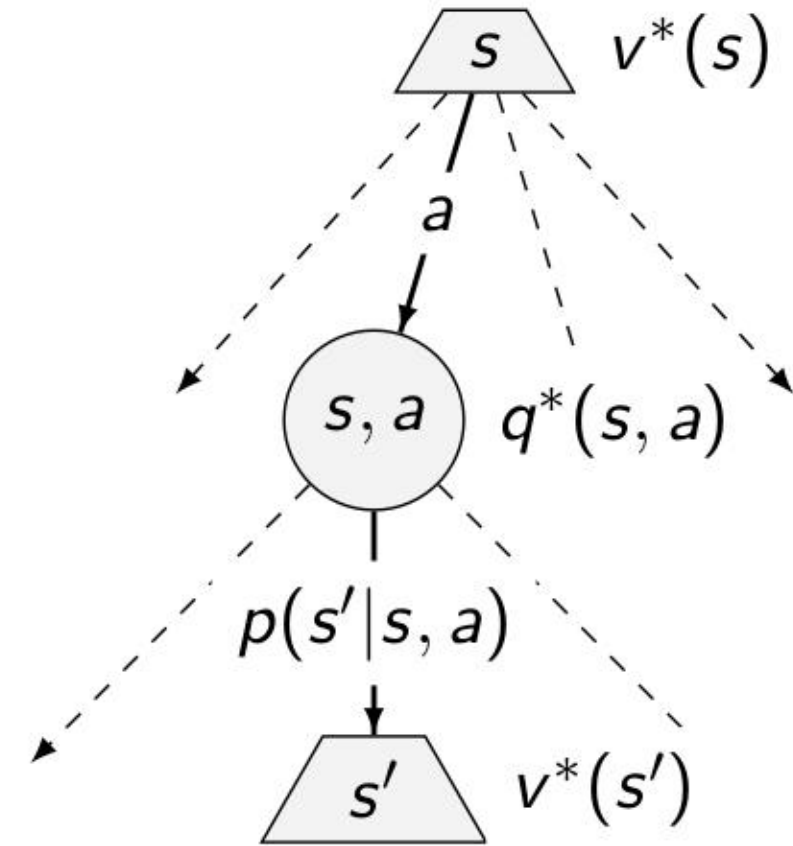
MDP search tree

The value of a q -state (s, a) :

$$q^*(s, a) = \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v^*(s')]$$

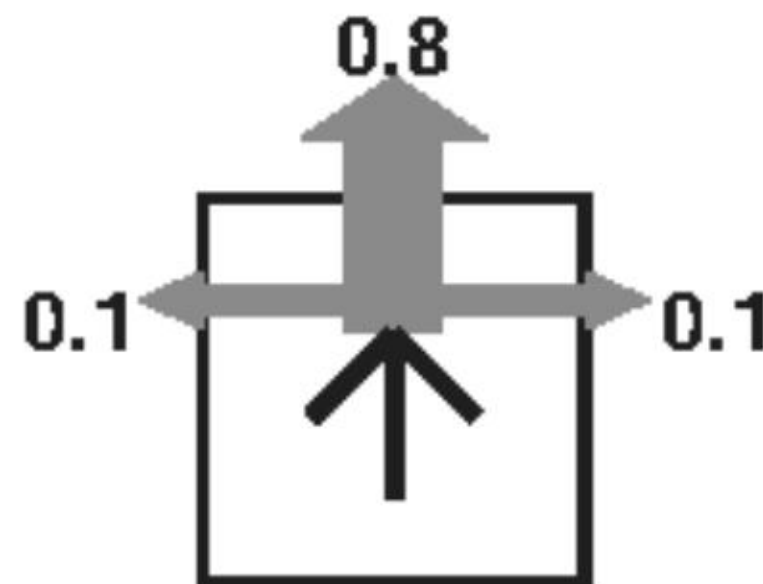
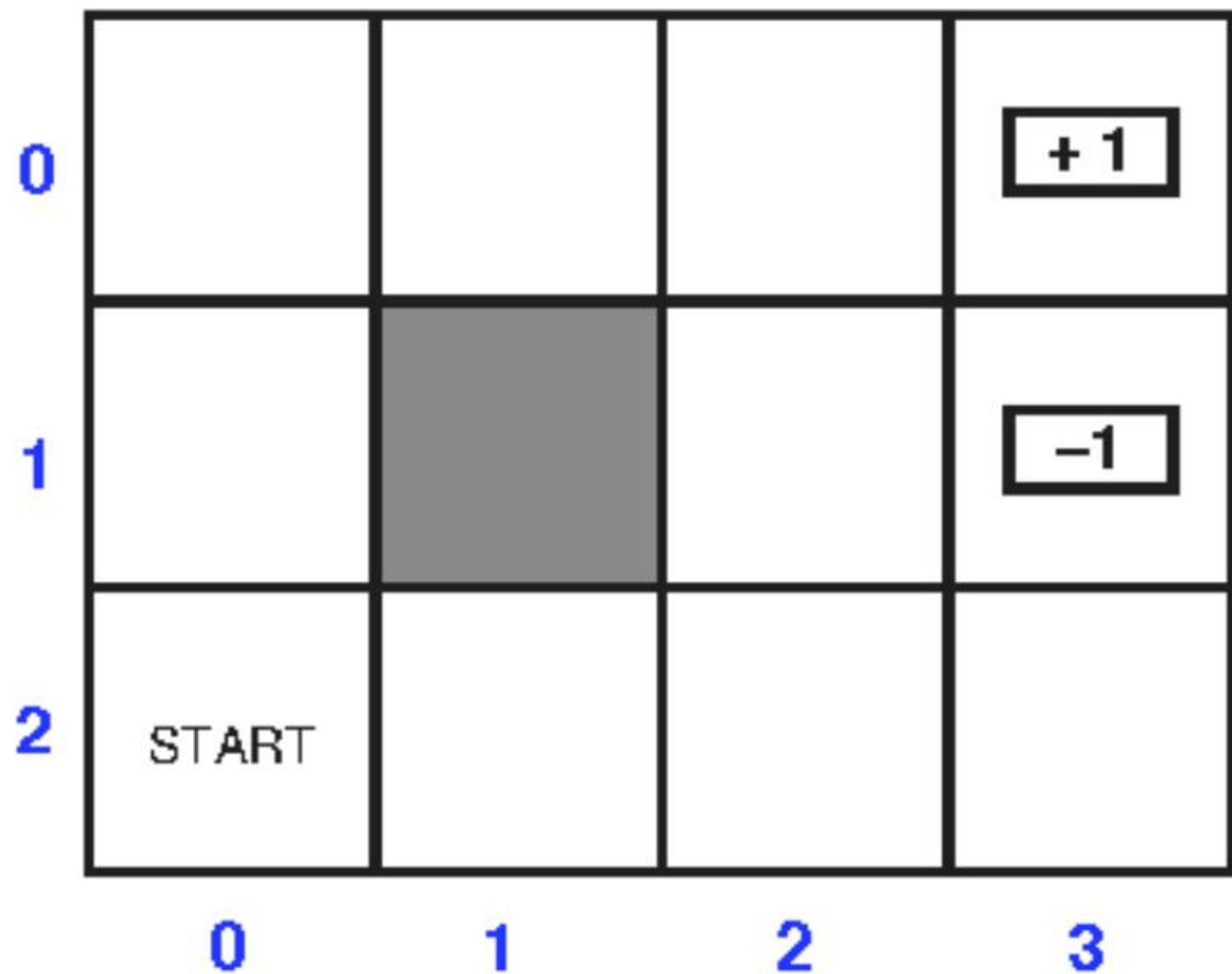
The value of a state s :

$$v^*(s) = \max_a q^*(s, a)$$



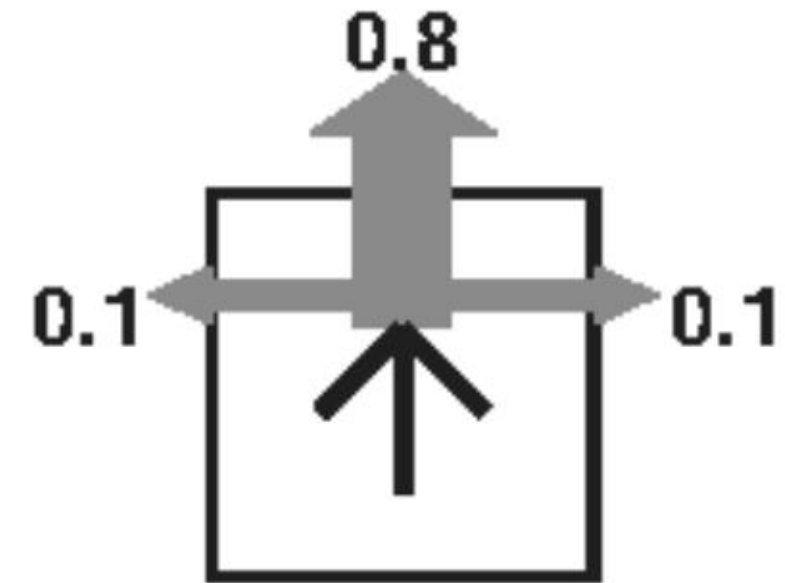
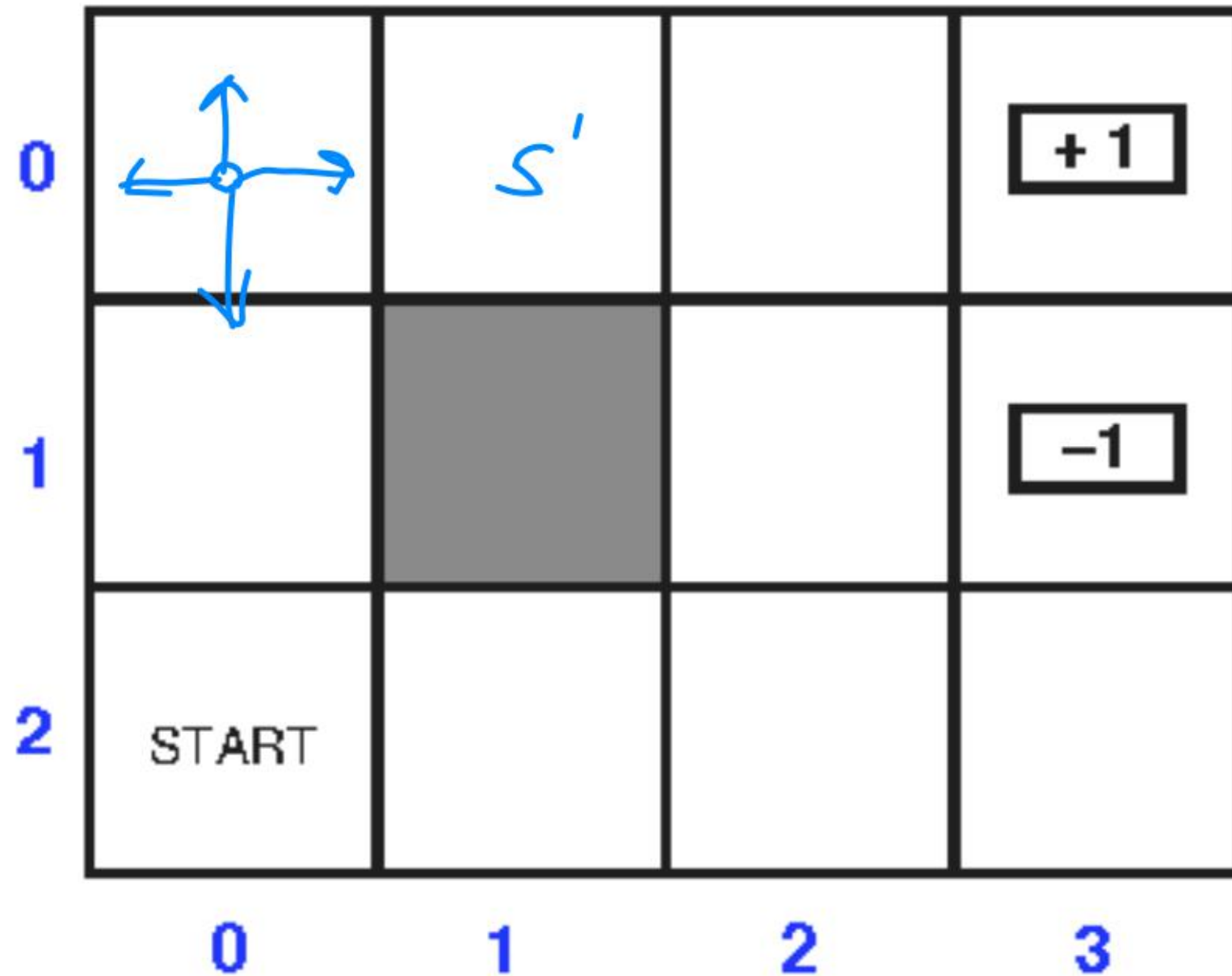
Bellman (optimality) equation

$$v^*(s) = \max_{a \in A(s)} \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v^*(s')]$$

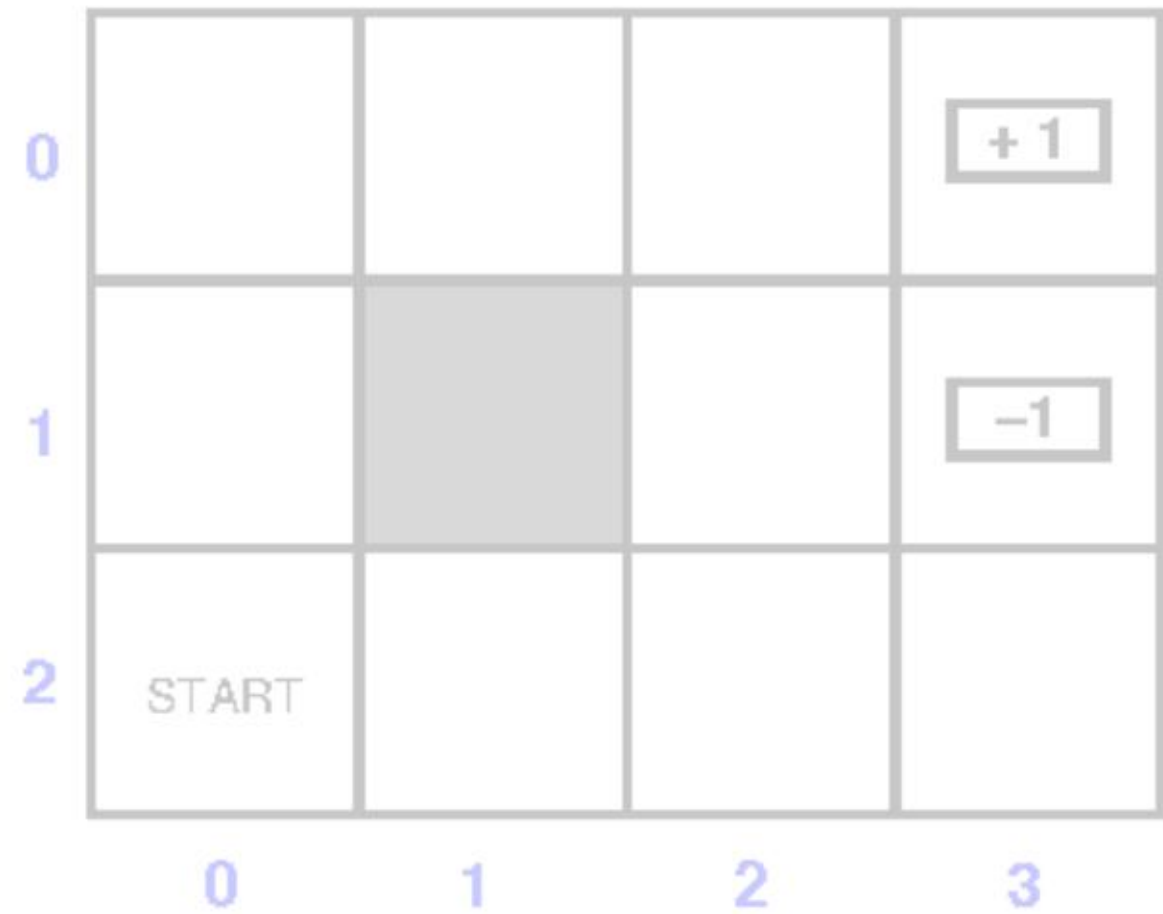


Bellman (optimality) equation

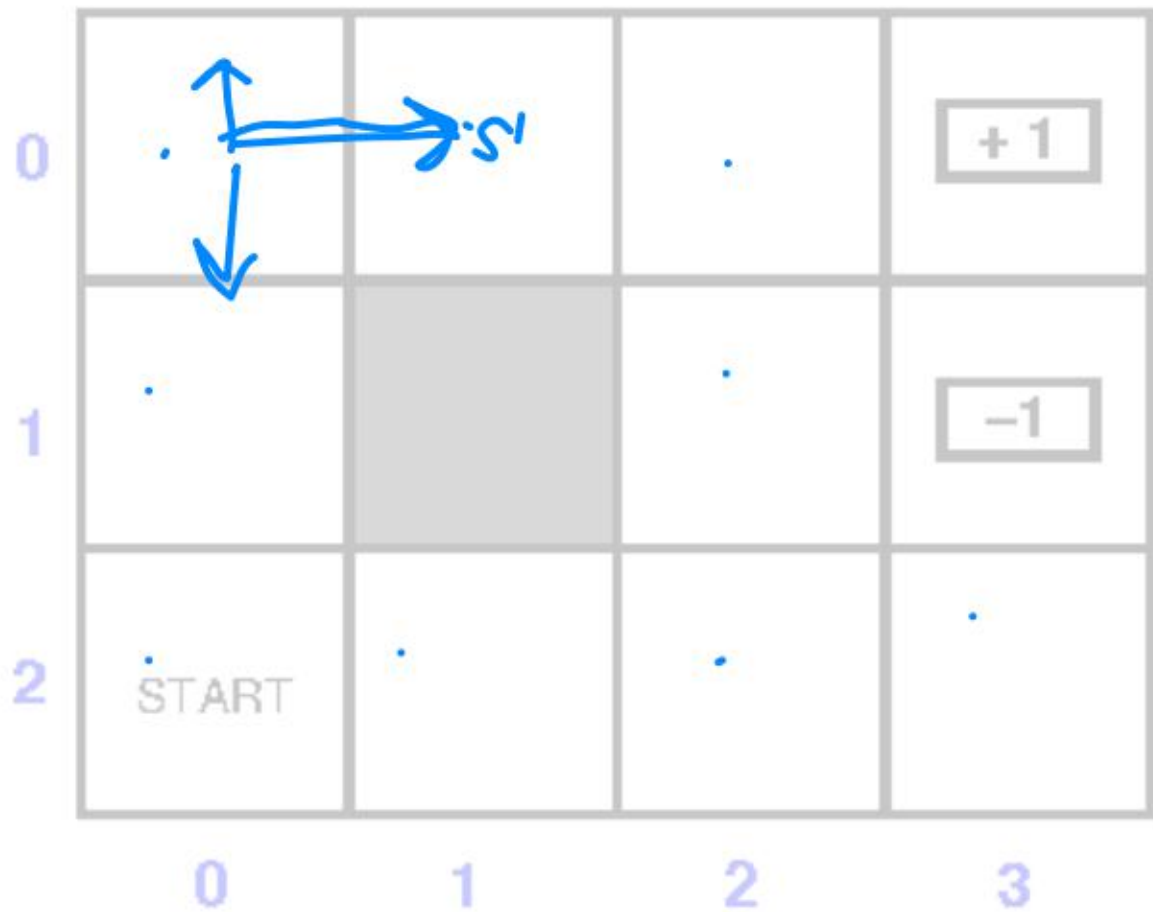
$$v^*(s) = \max_{a \in A(s)} \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v^*(s')]$$



$$v^*(s) = \max_{a \in A(s)} \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v^*(s')]$$



$$v^*(s) = \max_{a \in A(s)} \sum_{s'} p(s'|a, s) [r(s, a, s') + \gamma v^*(s')]$$



$$v(0,0) = a = E$$

$$s' 0,1: 0.8 [-0.04 + \gamma v(0,1)]$$

$$s' 0,0: 0.1 [-0.04 + \gamma v(0,0)]$$

$\gamma: \gamma$

Value iteration

- ▶ Start with arbitrary $V_0(s)$ (except for terminals)
- ▶ Compute Bellman update (one ply of expectimax from each state)

$$V_{k+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} p(s'|s,a) V_k(s')$$

- ▶ Repeat until convergence

The idea: Bellman update makes local consistency with the Bellmann equation. Everywhere locally consistent \Rightarrow globally optimal.

Value iteration

- ▶ Start with arbitrary $V_0(s)$ (except for terminals)
- ▶ Compute **Bellman update** (one ply of expectimax from each state)

$$V_{k+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} p(s'|s, a) V_k(s')$$

- ▶ Repeat until convergence

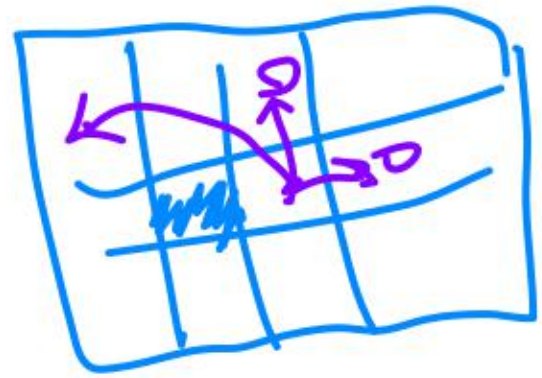
The idea: Bellman update makes local consistency with the Bellmann equation. Everywhere locally consistent \Rightarrow globally optimal.

Value iteration - Complexity of one estimation sweep

$$V_{k+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} p(s'|s, a) V_k(s')$$

- ▶ $O(AS)$
- ▶ $O(S^2)$
- ▶ $O(AS^2)$
- ▶ $O(A^2S^2)$

Value iteration - Complexity of one estimation sweep ;



$$* \underset{k}{V} \begin{bmatrix} 0 \\ 0 \\ \triangleright \\ b \end{bmatrix}$$

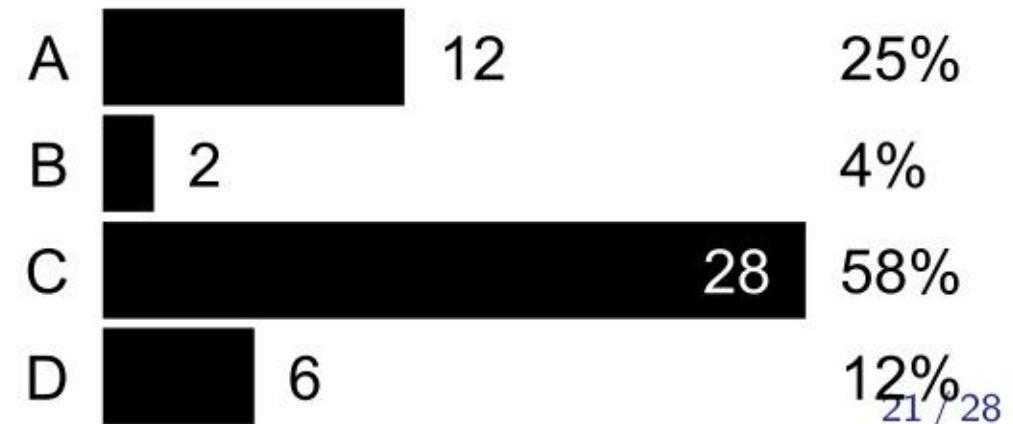
BV

$$\underset{k+1}{V'}$$

$$V_{k+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} p(s'|s, a) V_k(s')$$

SSA

- A ▶ $O(AS)$
- B ▶ $O(S^2)$
- C ▶ $O(AS^2)$
- D ▶ $O(A^2S^2)$



Convergence cont'd

$V_{k+1} \leftarrow BV_k \dots B$ as the Bellman update $V_{k+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} p(s'|s, a) V_k(s')$

$$\|BV_k - BV'_k\| \leq \gamma \|V_k - V'_k\|$$

$$\|BV_k - V_{\text{true}}\| \leq \gamma \|V_k - V_{\text{true}}\|$$

Rewards are bounded, at the beginning then Value error is

$$\|V_0 - V_{\text{true}}\| \leq \frac{2R_{\text{max}}}{1-\gamma}$$

We run N iterations and reduce the error by factor γ in each and want to stop the error is below ϵ :

$$\gamma^N 2R_{\text{max}} / (1 - \gamma) \leq \epsilon \text{ Taking logs, we find: } N \geq \frac{\log(2R_{\text{max}}/\epsilon(1-\gamma))}{\log(1/\gamma)}$$

To stop the iteration we want to find a bound relating the error to the size of *one* Bellman update for any given iteration.

We stop if

$$\|V_{k+1} - V_k\| \leq \frac{\epsilon(1-\gamma)}{\gamma}$$

then also: $\|V_{k+1} - V_{\text{true}}\| \leq \epsilon$ Proof on the next slide

Convergence cont'd

$\|V_{k+1} - V_{\text{true}}\| \leq \epsilon$ is the same as $\|V_{k+1} - V_{\infty}\| \leq \epsilon$

Assume $\|V_{k+1} - V_k\| = \text{err}$

In each of the following iteration steps we reduce the error by the factor γ (because $\|BV_k - V_{\text{true}}\| \leq \gamma\|V_k - V_{\text{true}}\|$). Till ∞ , the total sum of reduced errors is:

$$\text{total} = \gamma \text{err} + \gamma^2 \text{err} + \gamma^3 \text{err} + \gamma^4 \text{err} + \dots = \frac{\gamma \text{err}}{(1 - \gamma)}$$

We want to have $\text{total} < \epsilon$.

$$\frac{\gamma \text{err}}{(1 - \gamma)} < \epsilon$$

From it follows that

$$\text{err} < \frac{\epsilon(1 - \gamma)}{\gamma}$$

Hence we can stop if $\|V_{k+1} - V_k\| < \epsilon(1 - \gamma)/\gamma$

Convergence cont'd

$\|V_{k+1} - V_{\text{true}}\| \leq \epsilon$ is the same as $\|V_{k+1} - V_{\infty}\| \leq \epsilon$

Assume $\|V_{k+1} - V_k\| = \text{err}$

In each of the following iteration steps we reduce the error by the factor γ (because $\|BV_k - V_{\text{true}}\| \leq \gamma \|V_k - V_{\text{true}}\|$). Till ∞ , the total sum of reduced errors is:

$$\text{total} = \gamma \text{err} + \gamma^2 \text{err} + \gamma^3 \text{err} + \gamma^4 \text{err} + \dots = \frac{\gamma \text{err}}{(1 - \gamma)}$$

We want to have $\text{total} < \epsilon$.

$$\frac{\gamma \text{err}}{(1 - \gamma)} < \epsilon$$

From it follows that

$$\text{err} < \frac{\epsilon(1 - \gamma)}{\gamma}$$

Hence we can stop if $\|V_{k+1} - V_k\| < \epsilon(1 - \gamma)/\gamma$

Value iteration algorithm

function VALUE-ITERATION(env, ϵ) **returns:** state values V

input: env - MDP problem, ϵ

$V' \leftarrow 0$ in all states

repeat

$V \leftarrow V'$

$\delta \leftarrow 0$

for each state s in S do

$V'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) V(s')$

if $|V'[s] - V[s]| > \delta$ then $\delta \leftarrow |V'[s] - V[s]|$

end for

until $\delta < \epsilon(1 - \gamma)/\gamma$

end function

▷ iterate values until convergence

▷ keep the last known values

▷ reset the max difference

Value iteration algorithm

function VALUE-ITERATION(env, ϵ) **returns:** state values V

input: env - MDP problem, ϵ

$V' \leftarrow 0$ in all states

repeat

$V \leftarrow V'$

$\delta \leftarrow 0$

for each state s in S do

$V'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) V(s')$

if $|V'[s] - V[s]| > \delta$ then $\delta \leftarrow |V'[s] - V[s]|$

end for

until $\delta < \epsilon(1 - \gamma)/\gamma$

end function

- ▷ iterate values until convergence
- ▷ keep the last known values
- ▷ reset the max difference

Value iteration algorithm

function VALUE-ITERATION(env, ϵ) **returns:** state values V

input: env - MDP problem, ϵ

$V' \leftarrow 0$ in all states

repeat

— $V \leftarrow V'$

$\delta \leftarrow 0$

for each state s in S do

$V'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) V(s')$

if $|V'[s] - V[s]| > \delta$ then $\delta \leftarrow |V'[s] - V[s]|$

end for

until $\delta < \epsilon(1 - \gamma)/\gamma$

end function

- ▷ iterate values until convergence
- ▷ keep the last known values
- ▷ reset the max difference

Value iteration algorithm

function VALUE-ITERATION(env, ϵ) **returns:** state values V

input: env - MDP problem, ϵ

$V' \leftarrow 0$ in all states

repeat

$V \leftarrow V'$

$\delta \leftarrow 0$

for each state s **in** S **do**

$V'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) V(s')$

if $|V'[s] - V[s]| > \delta$ **then** $\delta \leftarrow |V'[s] - V[s]|$

end for

until $\delta < \epsilon(1 - \gamma)/\gamma$

end function

- ▷ iterate values until convergence
- ▷ keep the last known values
- ▷ reset the max difference