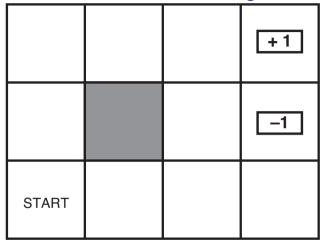
Sequential decisions under uncertainty Markov Decision Processes (MDP)

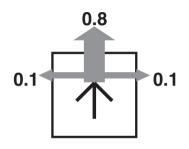
Tomáš Svoboda, Petr Pošík

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

March 19, 2025

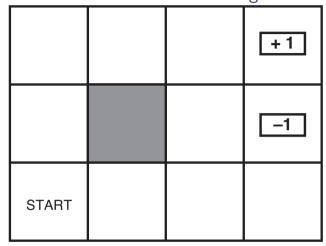
Unreliable actions in observable grid world

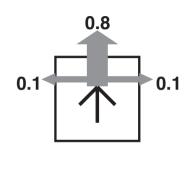




States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$ (Transition) Model $T(s, a, s') \equiv p(s'|s, a) =$ probability that a in s leads to s

Unreliable actions in observable grid world



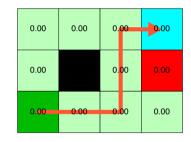


States $s \in \mathcal{S}$, actions $a \in \mathcal{A}$ (Transition) Model $T(s, a, s') \equiv p(s'|s, a) = \text{probability that } a \text{ in } s \text{ leads to } s'$

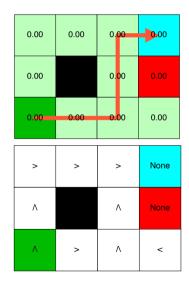
Unreliable (results of) actions



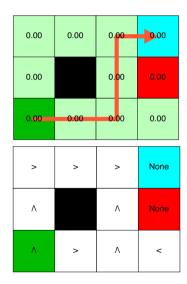
- ► In deterministic world: Plan sequence of actions from Start to Goal.
- ▶ MDPs, we need a policy $\pi: \mathcal{S} \to \mathcal{A}$
- An action for each possible state. Why each?
- What is the best policy?



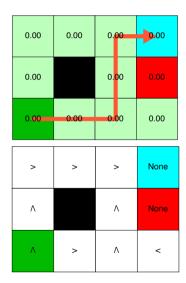
- ► In deterministic world: Plan sequence of actions from Start to Goal.
- ▶ MDPs, we need a policy $\pi: \mathcal{S} \to \mathcal{A}$.
- An action for each possible state. Why each?
- What is the best policy?



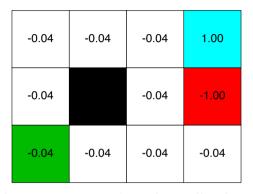
- ► In deterministic world: Plan sequence of actions from Start to Goal.
- ▶ MDPs, we need a policy $\pi: \mathcal{S} \to \mathcal{A}$.
- ► An action for each possible state. Why *each*?
- What is the best policy?



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



Rewards, Reward function

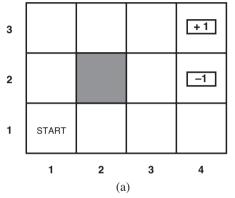


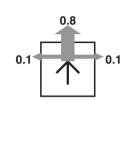
Reward : Robot/Agent takes an action a and it is **immediately** rewarded.

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

5 / 29

Markov Decision Processes (MDPs)



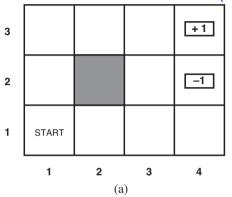


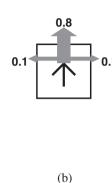
(b)

States $s \in S$, actions $a \in A$ Model $T(s, a, s') \equiv p(s'|s, a) = \text{probability that } a \text{ in } s$ Reward function r(s) (or r(s, a), r(s, a, s'))

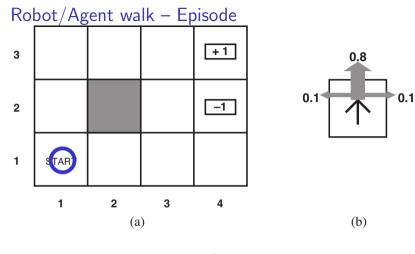
 $= \left\{ egin{array}{ll} -0.04 & ext{(small penalty) for nonterminal states} \ \pm 1 & ext{for terminal states} \end{array}
ight.$

Markov Decision Processes (MDPs)

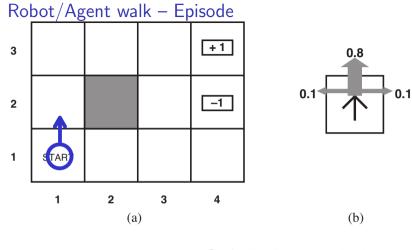




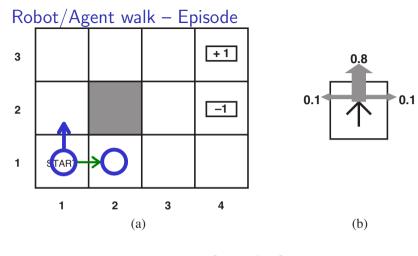
States
$$s \in \mathcal{S}$$
, actions $a \in \mathcal{A}$
Model $T(s, a, s') \equiv p(s'|s, a) = \text{probability that } a \text{ in } s \text{ leads to } s'$
Reward function $r(s)$ (or $r(s, a)$, $r(s, a, s')$)
$$= \begin{cases} -0.04 & \text{(small penalty) for nonterminal states} \\ \pm 1 & \text{for terminal states} \end{cases}$$



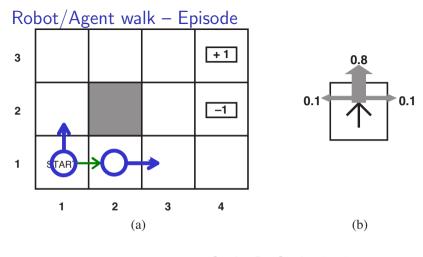
 $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2 \dots$



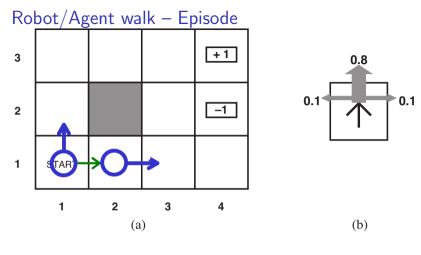
 $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2 \dots$



 $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2 \dots$



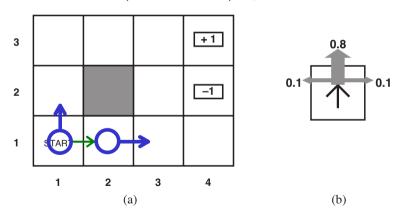
 $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2 \dots$



 $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2 \dots$

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.



Desired robot/agent behavior specified through rewards

- Before: shortest/cheapest path
- Solution found by search.
- Environment/problem is defined through the reward function
- Optimal policy is to be computed/learned.

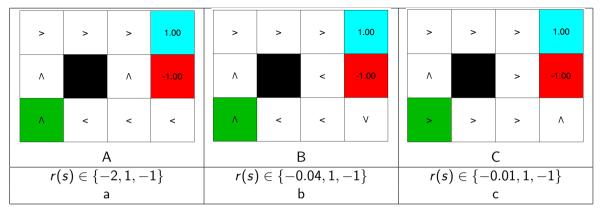
Desired robot/agent behavior specified through rewards

- ► Before: shortest/cheapest path
- Solution found by search.
- ► Environment/problem is defined through the reward function.
- ▶ Optimal policy is to be computed/learned.

Desired robot/agent behavior specified through rewards

- ► Before: shortest/cheapest path
- Solution found by search.
- ► Environment/problem is defined through the reward function.
- ▶ Optimal policy is to be computed/learned.

We come back to this in more detail when discussing RL.

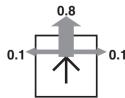


A: A-a, B-b, C-c

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

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

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



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

Typically, consider stationary preferences on reward sequences

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

Utility (h-history) $U_h([S_0,S_1,S_2,\ldots,])=R_1+R_2+R_3+\cdots$

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

Typically, consider stationary preferences on reward sequences:

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

Utility (h-history) $U_h([S_0,S_1,S_2,\ldots,])=R_1+R_2+R_3+\cdots$

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

Typically, consider stationary preferences on reward sequences:

$$[R,R_1,R_2,R_3,\ldots] \succ [R,R_1',R_2',R_3',\ldots] \Leftrightarrow [R_1,R_2,R_3,\ldots] \succ [R_1',R_2',R_3',\ldots]$$

```
If stationary preferences : Utility (h-history) U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \cdots
```

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

Typically, consider stationary preferences on reward sequences:

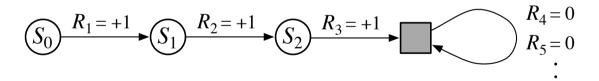
$$[R,R_1,R_2,R_3,\ldots] \succ [R,R_1',R_2',R_3',\ldots] \Leftrightarrow [R_1,R_2,R_3,\ldots] \succ [R_1',R_2',R_3',\ldots]$$

```
If stationary preferences : Utility (h-history) U_h([S_0, S_1, S_2, \dots,]) = R_1 + R_2 + R_3 + \cdots
```

Finite walk – Episode – and its Return (by introducing Terminal state)

- 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)

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



Horizon too far, infinite - Discount rewards

Problem: Infinite lifetime ⇒ additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- Absorbing (terminal) state. (sooner or later walk ends here)
- ightharpoonup Discounted return , $\gamma < 1, R_t \le R_{\text{max}}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{1} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

Horizon too far, infinite – Discount rewards

Problem: Infinite lifetime ⇒ additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- Absorbing (terminal) state. (sooner or later walk ends here
- ightharpoonup Discounted return , $\gamma < 1, R_t \le R_{\max}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{1} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

Horizon too far, infinite – Discount rewards

Problem: Infinite lifetime ⇒ additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- Absorbing (terminal) state. (sooner or later walk ends here)
- ▶ Discounted return , $\gamma < 1, R_t \le R_{\text{max}}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{1} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

Horizon too far, infinite - Discount rewards

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.
- Absorbing (terminal) state. (sooner or later walk ends here)
- ▶ Discounted return , $\gamma < 1, R_t \le R_{\mathsf{max}}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{1} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

Horizon too far, infinite - Discount rewards

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.
- Absorbing (terminal) state. (sooner or later walk ends here)
- lacktriangle Discounted return , $\gamma < 1, R_t \leq R_{\sf max}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{4} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

Horizon too far, infinite – Discount rewards

Problem: Infinite lifetime ⇒ additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- ► Absorbing (terminal) state. (sooner or later walk ends here)
- ▶ Discounted return , $\gamma < 1, R_t \le R_{\mathsf{max}}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{1} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

Horizon too far, infinite - Discount rewards

Problem: Infinite lifetime ⇒ additive utilities are infinite.

- ▶ Finite horizon: termination at a fixed time \Rightarrow nonstationary policy, $\pi(s)$ depends on the time left.
- Absorbing (terminal) state. (sooner or later walk ends here)
- ▶ Discounted return , $\gamma < 1, R_t \le R_{\sf max}$

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

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma^{1} R_{t+3} + \gamma^{2} R_{t+4} + \cdots)$$

$$= R_{t+1} + \gamma G_{t+1}$$

MDPs recap

Markov decision processes (MDPs):

- \triangleright Set of states \mathcal{S}
- \triangleright Set of actions \mathcal{A}
- ▶ Transitions p(s'|s, a) or T(s, a, s')
- ▶ Reward function r(s, a, s'); and discount γ
- Alternative to last two: p(s', r|s, a).

MDP quantities:

- lacktriangle (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):

- \triangleright Set of states \mathcal{S}
- \triangleright Set of actions \mathcal{A}
- ► Transitions p(s'|s, a) or T(s, a, s')
- ▶ Reward function r(s, a, s'); and discount γ
- Alternative to last two: p(s', r|s, a).

MDP quantities:

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

Expected Return of a policy π

- ightharpoonup Executing policy $\pi \to \text{sequence of states (and rewards)}.$
- Utility of a state sequence.
- But actions are unreliable environment is stochastic
- ightharpoonup Expected return of a policy π .

Starting at time t, i.e. S_t ,

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

Expected Return of a policy π

- ightharpoonup Executing policy $\pi \to \text{sequence of states (and rewards)}.$
- Utility of a state sequence.
- But actions are unreliable environment is stochastic.
- ightharpoonup Expected return of a policy π .

Starting at time t, i.e. S_t ,

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

Expected Return of a policy π

- ightharpoonup Executing policy $\pi \to \text{sequence of states (and rewards)}.$
- 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 ,

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

(State) Value functions given policy π

Expected return from that state (state, action)

Value function

$$v^\pi(s) \stackrel{def}{=} \mathsf{E}^\pi \left[\mathsf{G}_t \mid \mathsf{S}_t = \mathsf{s}
ight] = \mathsf{E}^\pi \left[\sum_{k=0}^\infty \gamma^k \mathsf{R}_{t+k+1} \; \middle| \; \mathsf{S}_t = \mathsf{s}
ight]$$

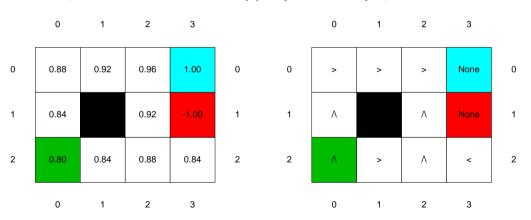
Action-value function (q-function)

$$q^{\pi}(s,a) \stackrel{def}{=} \mathsf{E}^{\pi}\left[\mathsf{G}_t \mid \mathsf{S}_t = \mathsf{s}, \mathsf{A}_t = \mathsf{a}\right] = \mathsf{E}^{\pi}\left[\sum_{k=0}^{\infty} \gamma^k \mathsf{R}_{t+k+1} \mid \mathsf{S}_t = \mathsf{s}, \mathsf{A}_t = \mathsf{a}\right]$$

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

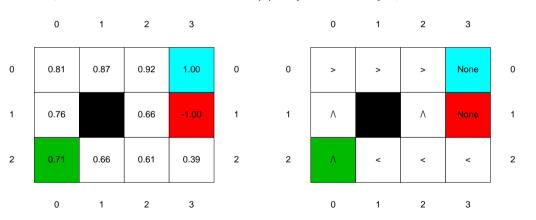
 $v^*(s) = \text{expected (discounted)}$ sum of rewards (until termination) assuming optimal actions.

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



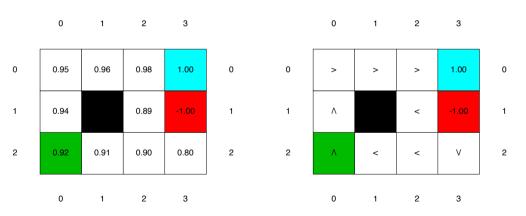
 $v^*(s) = \text{expected (discounted)}$ sum of rewards (until termination) assuming optimal actions.

Example 2, Robot *non-deterministic*: $r(s) = \{-0.04, 1, -1\}, \gamma = 0.9999999, \epsilon = 0.03$

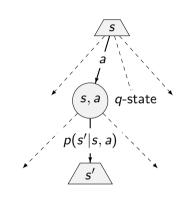


 $v^*(s) = \text{expected (discounted)}$ sum of rewards (until termination) assuming optimal actions.

Example 3, Robot *non-deterministic*: $r(s) = \{-0.01, 1, -1\}, \gamma = 0.9999999, \epsilon = 0.03$



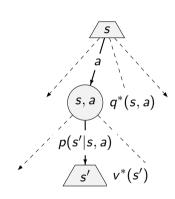
MDP search tree (Expectimax search)



MDP search tree (Expectimax search)

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

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



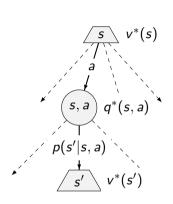
MDP search tree (Expectimax search)

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

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

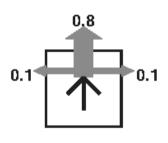
The value of a state s:

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



Bellman (optimality) equation

$$v^{*}(s) = \max_{a \in A(s)} \sum_{s'} p(s'|a, s) \left[r(s, a, s') + \gamma v^{*}(s') \right]$$
0
1
2
START
0.1



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

- 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.

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

- 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.

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

- 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.

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

- ▶ 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

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

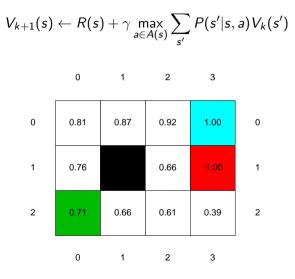
A: O(AS)

B: $O(S^2)$

 $C: O(AS^2)$

D: $O(A^2S^2)$

Value iteration demo



Convergence

$$egin{aligned} V_{k+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) V_k(s') \ & \gamma < 1 \ & -R_{\mathsf{max}} \leq R(s) \leq R_{\mathsf{max}} \end{aligned}$$

Max norm

$$\|v\|_{\infty} = \max_{s} |v(s)|$$
 $U([s_0, s_1, s_2, \dots, s_{\infty}]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq \frac{R_{\max}}{1 - \gamma}$

Convergence

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

$$\gamma < 1$$

$$-R_{\mathsf{max}} \leq R(s) \leq R_{\mathsf{max}}$$

Max norm:

$$\|V\|_{\infty} = \max_{s} |V(s)|$$
 $U([s_0, s_1, s_2, \dots, s_{\infty}]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq \frac{R_{\mathsf{max}}}{1-\gamma}$

Convergence cont'd

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

$$||BV_k - BV_k'||_{\infty} \le \gamma ||V_k - V_k'||_{\infty}$$
$$||BV_k - V_{\text{true}}||_{\infty} < \gamma ||V_k - V_{\text{true}}||_{\infty}$$

Rewards are bounded, at the beginning then Value error is

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

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

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

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

If we stop when

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

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

Convergence cont'd

$$\|V_{k+1} - V_{\text{true}}\|_{\infty} \le \epsilon$$
 is the same as $\|V_{k+1} - V_{\infty}\|_{\infty} \le \epsilon$
Assume $\|V_{k+1} - V_k\|_{\infty} = \text{err}$

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

total =
$$\gamma$$
err + γ^2 err + γ^3 err + γ^4 err + \cdots = $\frac{\gamma$ err γ^4 err + γ^4 err +

We want to have total $< \epsilon$.

$$rac{\gamma \mathsf{err}}{(1-\gamma)} < \epsilon$$

From it follows that

$$\operatorname{\mathsf{err}} < rac{\epsilon(1-\gamma)}{\gamma}$$

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

```
function VALUE-ITERATION(env,\epsilon) returns: state values V
    input: env - MDP problem. \epsilon
    V' \leftarrow 0 in all states
```

```
function VALUE-ITERATION(env,\epsilon) returns: state values V
   input: env - MDP problem. \epsilon
   V' \leftarrow 0 in all states
                                                    repeat
```

```
function VALUE-ITERATION(env,\epsilon) returns: state values V
   input: env - MDP problem, \epsilon
   V' \leftarrow 0 in all states
                                                 repeat
      V \leftarrow V'
                                           \delta \leftarrow 0
                                                       > reset the max difference
```

```
function VALUE-ITERATION(env,\epsilon) returns: state values V
   input: env - MDP problem, \epsilon
    V' \leftarrow 0 in all states
                                                             repeat
       V \leftarrow V'
                                                     \delta \leftarrow 0
                                                                     > reset the max difference
       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]|
```

```
function VALUE-ITERATION(env,\epsilon) returns: state values V
   input: env - MDP problem, \epsilon
    V' \leftarrow 0 in all states
                                                                repeat
        V \leftarrow V'
                                                        \delta \leftarrow 0
                                                                        > reset the max difference
        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]|
   until \delta < \epsilon (1 - \gamma)/\gamma
```

Sync vs. async Value iteration

```
function VALUE-ITERATION(env,\epsilon) returns: state values V
   input: env - MDP problem, \epsilon
    V' \leftarrow 0 in all states
                                                                repeat
        V = V'
                                                              \delta \leftarrow 0
                                                                        > reset the max difference
        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]|
   until \delta < \epsilon (1 - \gamma)/\gamma
```

What we have learned

- Uncertain outcome of an action
- Optimal policy (strategy, sequence of decisions) maximizes expected return (utility, sum of rewards)
- ► (State) Value function given policy
- ▶ Value iteration method through local (optimal) updated to global optimality

References

Some figures from [1] (chapter 17) but notation slightly changed in order to adapt notation from [2] (chapters 3, 4) which will help us in the Reinforcement Learning part of the course. Note that the book [2] is available on-line.

[1] Stuart Russell and Peter Norvig.

Artificial Intelligence: A Modern Approach.

Prentice Hall, 3rd edition, 2010.

http://aima.cs.berkeley.edu/.

[2] Richard S. Sutton and Andrew G. Barto. Reinforcement Learning; an Introduction.

MIT Press, 2nd edition, 2018.

http://www.incompleteideas.net/book/the-book-2nd.html.