**Học phần RBE3043: Các thuật toán thích nghi**

# Buổi 4: Quá trình ra quyết định Markov hữu hạn (cont.)

**Giảng viên: TS. Nguyễn Thế Hoàng Anh**

*Hà Nội, ngày 11 tháng 10 năm 2023*
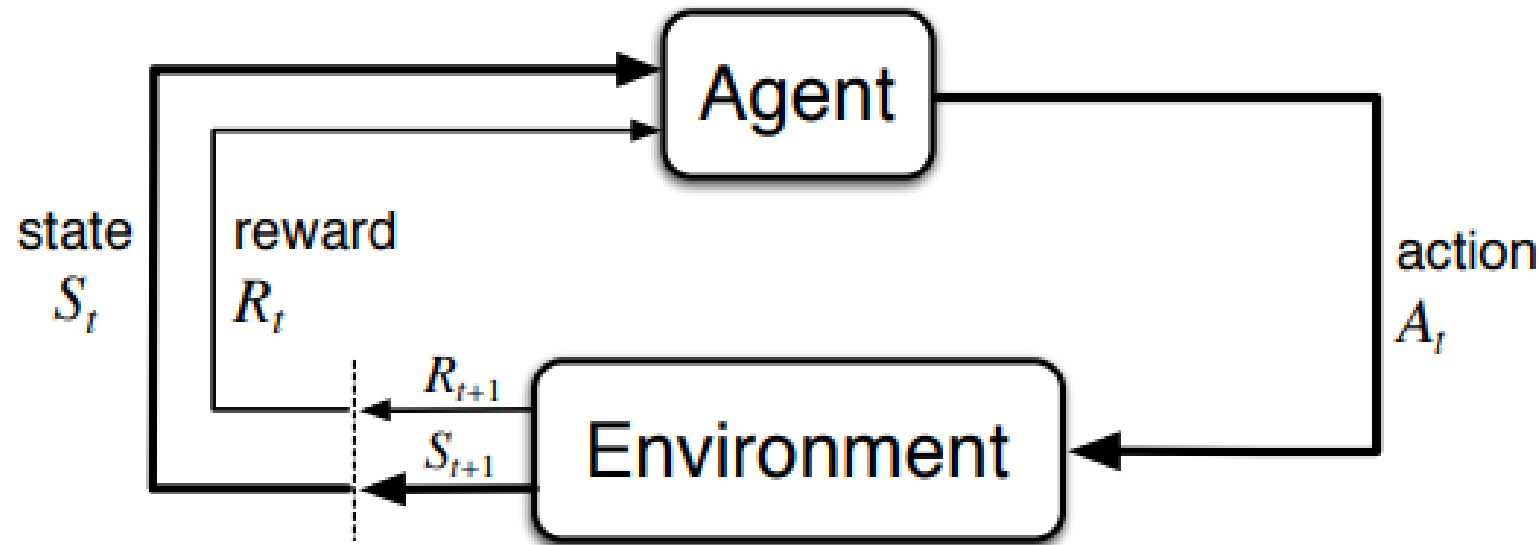
# The Agent-Environment Interface



Figure 3.1: The agent–environment interaction in a Markov decision process.

MDPs are meant to be a straightforward framing of the problem of learning from interaction to achieve a goal
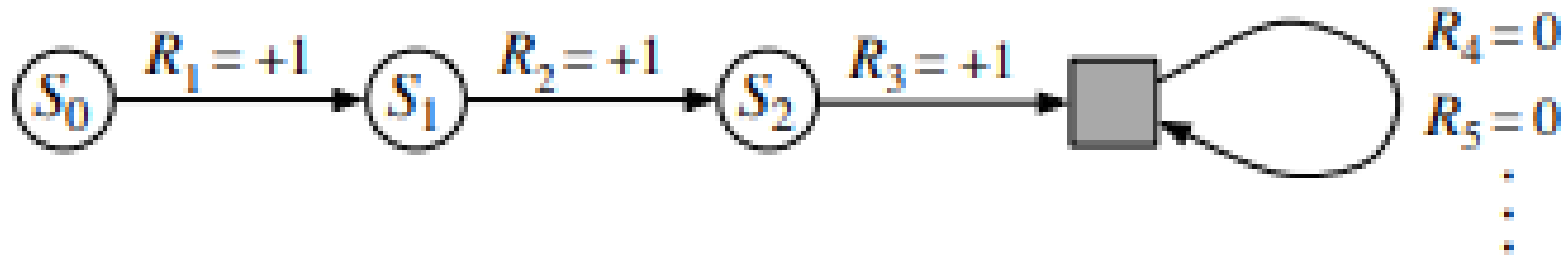
# MDP Tuple: <S, A, P, R, $\gamma$>

- *S: State of the agent on the grid (4,3)*
  - Note that cell denoted by (x,y)
- *A: Actions of the agent, i.e., N, E, S, W*

- *P: Transition function*
  - Table P(s' | s, a), prob of s' given action "a" in state "s"
  - E.g., P( (4,3) | (3,3), N) = 0.1
  - E.g., P((3, 2) | (3,3), N) = 0.8
  - (Robot movement, uncertainty of another agent's actions,…)

- *R: Reward  (more comments on the reward function later)*
  - *R( (3, 3), N) = -1/25*
  - R (4,1) = +1
- *$\gamma$: Discounted* factor

# Episodic and Continuing

- **Episodic tasks**: Reinforcement learning tasks that agent interact with environment within a sequence of separate episodes ~ have a terminal state
  - ✓Playing Go, Dota, Atari 2600
  - ✓Bioreactor
  - ✓Pic-and-place robot
  - ✓UAV
  - ✓Self-driving car
  - ✓Incheon airport robot

- **Continuing tasks**: others (not break down into episodes) ~ continue indefinitely ~ no terminal state
  - ✓Cart-pole
  - ✓Cycling robot
  - ✓Studying

# Episodic tasks

Episodic tasks have the followed type of state transition diagram



Absorbing state: transition to itself with zero reward → R4 - R5 - …

Return of episodic MDP

$$G_t \doteq \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k,$$

including the possibility that $T = \infty$ or $\gamma = 1$ (but not both).

If reward is constant +1 and the MDP is infinite

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}.$$

# Policies and Value functions

$v_\pi$ the *state-value function for policy* $\pi$

The *value* of a state $s$ under a policy $\pi$, denoted $v_\pi(s)$, is the expected return when starting in $s$ and following $\pi$ thereafter. For MDPs, we can define $v_\pi$ formally by

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right], \quad \text{for all } s \in \mathcal{S}, \tag{3.12}$$

$q_\pi$ the *action-value function for policy* $\pi$.

the value of taking action $a$ in state $s$ under a policy $\pi$, denoted $q_\pi(s,a)$, as the expected return starting from $s$, taking the action $a$, and thereafter following policy $\pi$:

$$q_\pi(s,a) \doteq \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a\right]. \tag{3.13}$$

These value functions measure how good it is for the agent to be in a given state or how good it is to perform a given action in a given state

# Bellman equation

Relationship between the value of a state and the values of its successor states

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s]$$
$$= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \qquad \text{(by (3.9))}$$
$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) \Big[ r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s'] \Big]$$
$$= \sum_a \pi(a|s) \sum_{s',r} p(s', r | s, a) \Big[ r + \gamma v_\pi(s') \Big], \quad \text{for all } s \in \mathcal{S}, \qquad (3.14)$$

agent is following policy $\pi$ at time $t$, then $\pi(a|s)$ is the probability that $A_t = a$ if $S_t = s$.

Meaning of Bellman equation: value of the start state must equal the (discounted) value of the expected next state plus the reward expected in the future

# Examples

**Example 3.5: Gridworld** Figure 3.2 (left) shows a rectangular gridworld representation of a simple finite MDP. The cells of the grid correspond to the states of the environment. At each cell, four actions are possible: north, south, east, and west, which deterministically cause the agent to move one cell in the respective direction on the grid. Actions that would take the agent off the grid leave its location unchanged, but also result in a reward of −1. Other actions result in a reward of 0, except those that move the agent out of the special states A and B. From state A, all four actions yield a reward of +10 and take the agent to A′. From state B, all actions yield a reward of +5 and take the agent to B′.

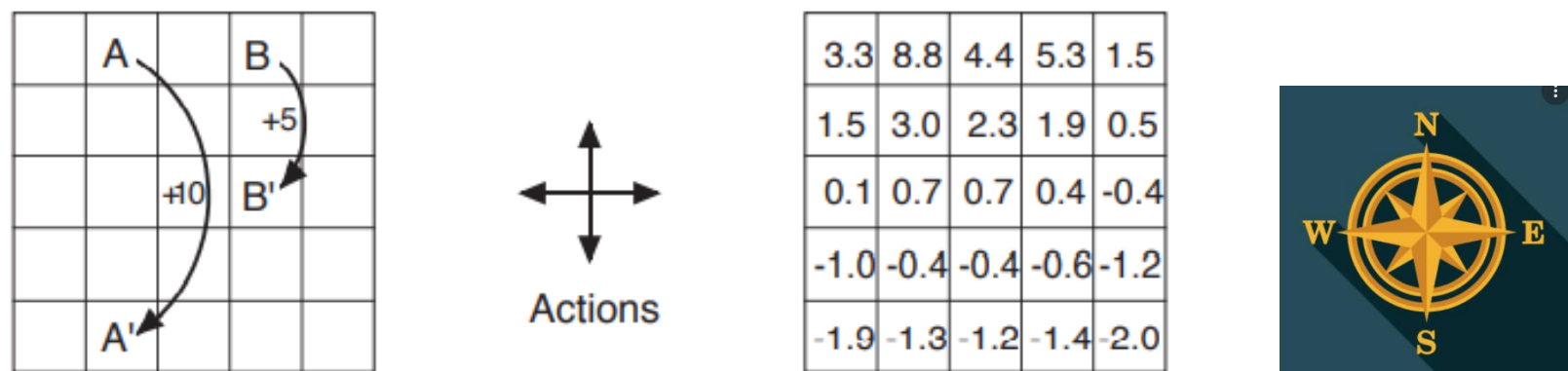| | | | | |
|---|---|---|---|---|
| 3.3 | 8.8 | 4.4 | 5.3 | 1.5 |
| 1.5 | 3.0 | 2.3 | 1.9 | 0.5 |
| 0.1 | 0.7 | 0.7 | 0.4 | -0.4 |
| -1.0 | -0.4 | -0.4 | -0.6 | -1.2 |
| -1.9 | -1.3 | -1.2 | -1.4 | -2.0 |

Actions

Figure 3.2: Gridworld example: exceptional reward dynamics (left) and state-value function for the equiprobable random policy (right).

Suppose the agent selects all four actions with equal probability in all states. Figure 3.2 (right) shows the value function, $v_\pi$, for this policy, for the discounted reward case with $\gamma = 0.9$. This value function was computed by solving the system of linear equations (3.14).

# Exercise 3.12

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| 10 | 9 | 8 | 7 | 6 |
| 11 | 12 | 13 | 14 | 15 |
| 20 | 19 | 18 | 17 | 16 |
| 21 | 22 | 23 | 24 | 25 |



| | | | | |
|---|---|---|---|---|
| | | 2.3 | | |
| | 0.7 | ? | 0.4 | |
| | | -0.4 | | |
| | | | | |

$v_\pi(s' = 12) = 0.7$          $v_\pi(s' = 18) = -0.4$

$v_\pi(s' = 14) = 0.4$          $v_\pi(s' = 8) = 2.3$

$v_\pi(s = 13) = ?$

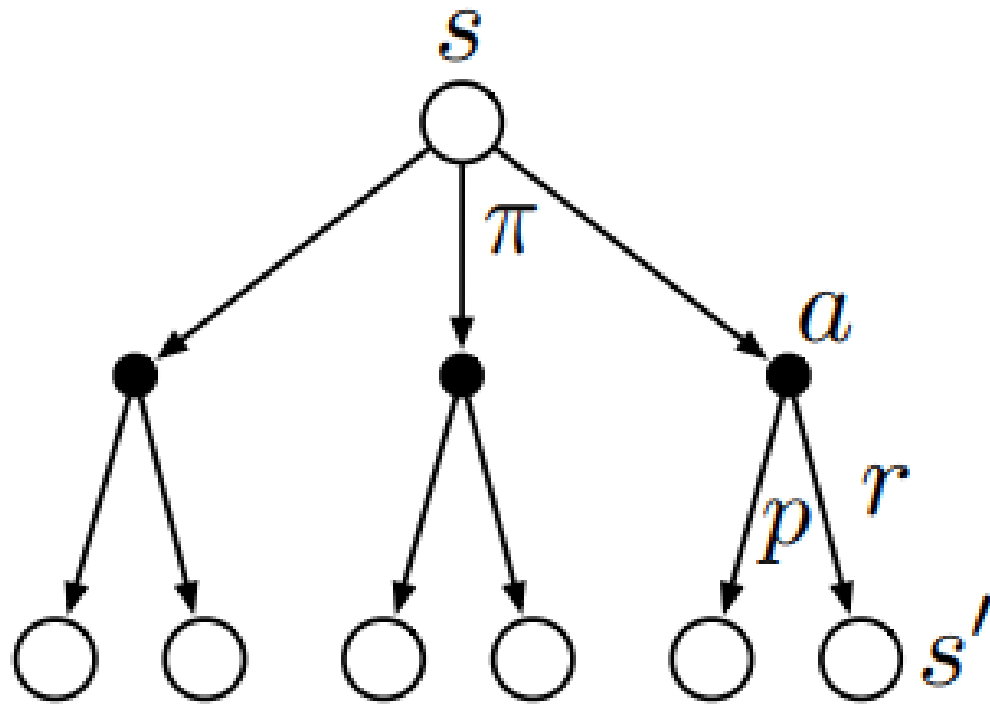$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\,[r + \gamma\,v_\pi(s')]$$

$$= \sum_a \pi(a|s) \sum_{s'} p(s',r=0|s,a)\,[r=0 + \gamma\,v_\pi(s')]$$

$$= \sum_a \pi(a|s) \sum_{s'} \gamma\,v_\pi(s')$$

$$= \sum_a \pi(a|s)\,\gamma \sum_{s'} v_\pi(s')$$

$$= \sum_a \pi(a|s) \times 0.9 \times \left[ v_\pi(s'=12) + v_\pi(s'=14) + v_\pi(s'=13) + v_\pi(s'=23) \right]$$

# Backup diagram



Backup diagram for $v_\pi$

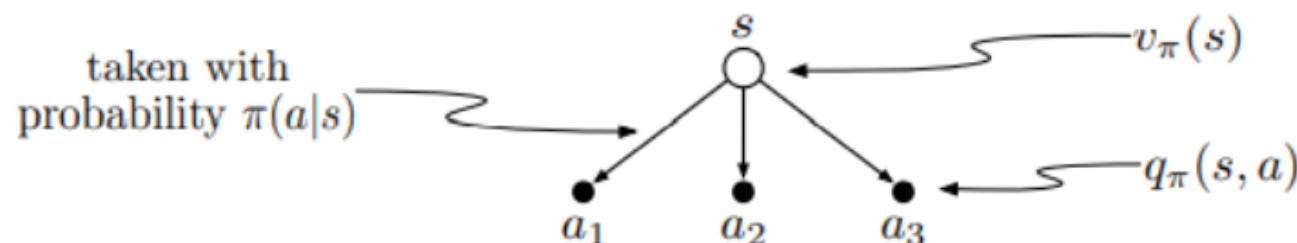Backup diagram: shows relationship between two successive states
Open circle: a state
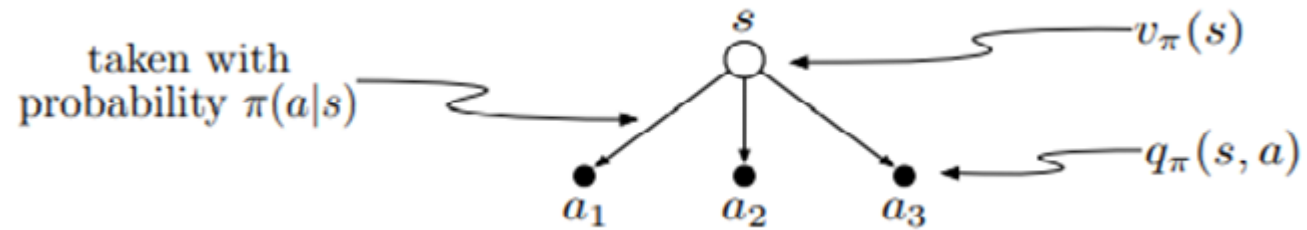Solid circle: a state-action pair
Policy: π

Backup (update) operations transfer value information back to a state (or a state-action pair) from its successor states (or state-action pairs)

*Exercise 3.16* The value of a state depends on the values of the actions possible in that state and on how likely each action is to be taken under the current policy. We can think of this in terms of a small backup diagram rooted at the state and considering each possible action:



taken with probability $\pi(a|s)$

$s$

$v_\pi(s)$

$q_\pi(s, a)$

$a_1 \quad a_2 \quad a_3$

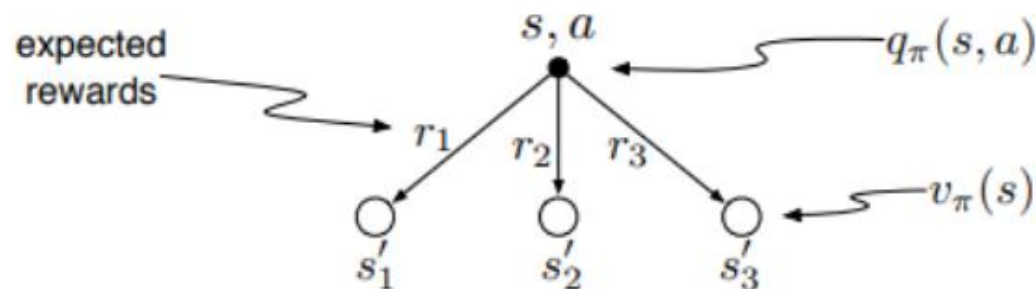Give the equation corresponding to this intuition and diagram for the value at the root node $v_\pi(s)$ in

$$v_\pi(s) = ?$$

$$v_\pi(s) = \pi(a_1|s)q_\pi(s|a_1) + \pi(a_2|s)q_\pi(s|a_2) + \pi(a_3|s)q_\pi(s|a_3)$$

$$v_\pi(s) = \sum_{a \in A} \pi(a|s)q_\pi(s,a)$$

*Exercise 3.17* The value of an action, $q_\pi(s, a)$, depends on the expected next reward and the expected sum of the remaining rewards. Again we can think of this in terms of a small backup diagram, this one rooted at an action (state–action pair) and branching to the possible next states:



Give the equation corresponding to this intuition and diagram for the action value, $q_\pi(s, a)$, in terms of the expected next reward, $R_{t+1}$, and the expected next state value, $v_\pi(S_{t+1})$, given that $S_t = s$ and $A_t = a$. This equation should include an expectation but *not* one conditioned on following the policy. Then give a second equation, writing out the expected value explicitly in terms of $p(s', r | s, a)$ defined by (3.2), such that no expected value notation appears in the equation. □

$$q_\pi(s, a) = ?$$

The action value $q(s,a)$ depends on expected next reward and expected sum of remaining rewards. Writing in terms of expectation we get

$$q_\pi(s,a) = \mathbb{E}_\pi\left[G_t | S_t = s, A_t = a\right]$$
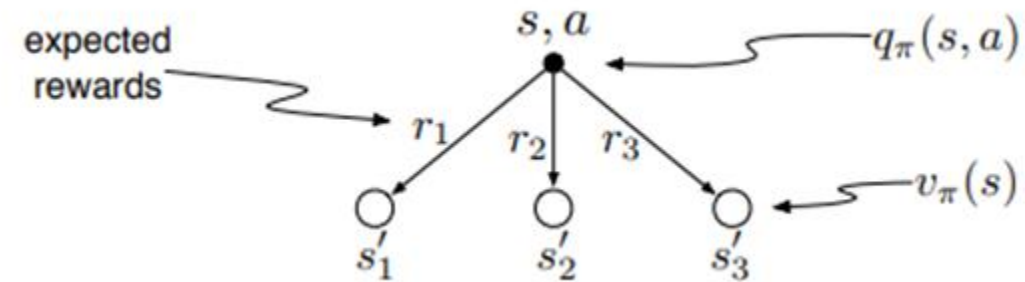
For all branches we get values as follows

$$q_\pi^1(s,a) = p(s_1, r_1 | s, a)\left[r_1 + \gamma v_\pi(s_1)\right]$$

$$q_\pi^2(s,a) = p(s_2, r_2 | s, a)\left[r_2 + \gamma v_\pi(s_2)\right]$$
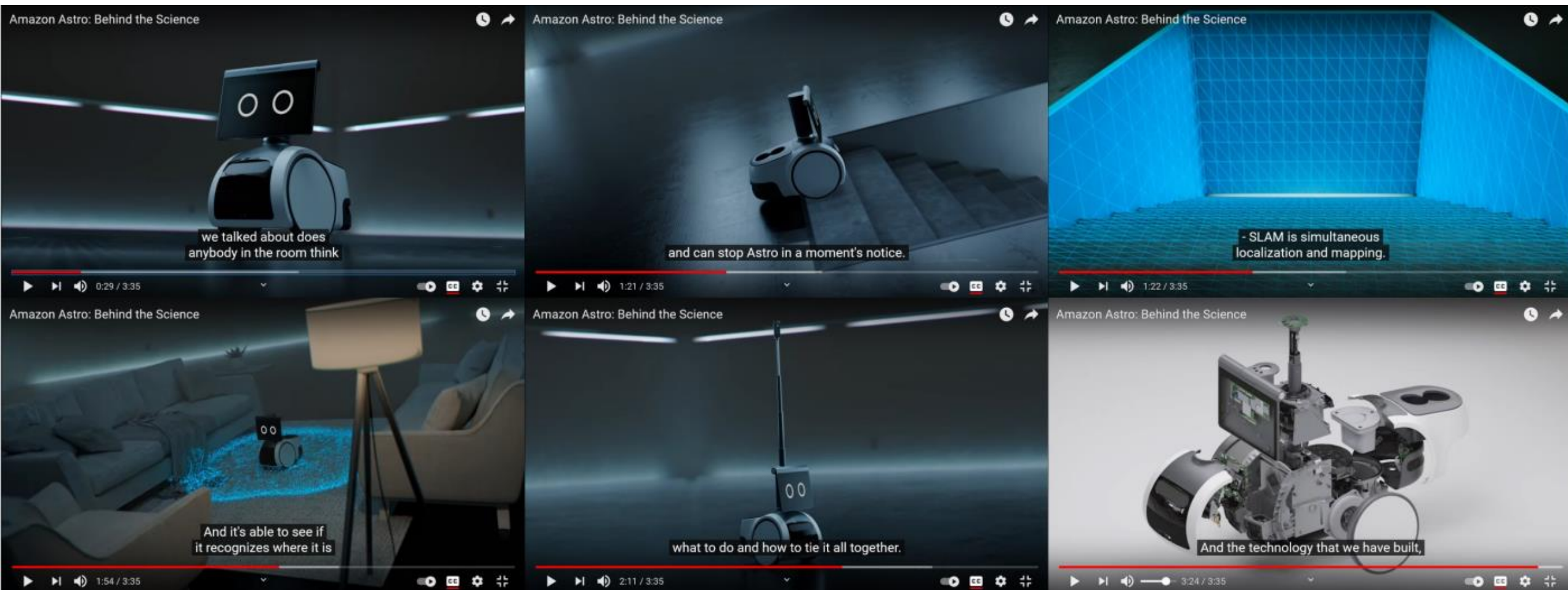
$$q_\pi^3(s,a) = p(s_3, r_3 | s, a)\left[r_3 + \gamma v_\pi(s_3)\right]$$

The final value will be sum of the three

$$q_\pi(s,a) = q_\pi^1(s,a) + q_\pi^2(s,a) + q_\pi^3(s,a)$$

$$= \sum_{s',r} p(s', r | s, a)\left[r + \gamma v_\pi(s')\right]$$

# Amazon Astro



Amazon Astro      https://www.youtube.com/watch?v=sj1t3msy8dc

Behind the Science   https://www.youtube.com/watch?v=Zy4If8-Wth4

# Amazon Astro: Technologies

2 camera (front) + 2 additional cameras, follow people

Computer vision

ROBOTICS

2 wheels, autonomous driving, SLAM – Simultaneous localization and mapping (reinforcement learning?), automatic charging

Robot understanding & communication

Information Retrieval
Sentiment Analysis
Information Extraction
Machine Translation
**Natural Language Processing (NLP)**
Question Answering
Google

**See**
**Infer**
**Hear**
**Understanding**
**Moving**

Speech recognition and processing: window break, sound alarm

Image processing: object detection, smoke detection

Face++ = new face recognition ();

Cloud Computing

video streaming

# Optimal Policies and Optimal Value Functions

Better policy: Expected return of policy $\pi$ is greater than or equal to expected return of policy $\pi'$ for all states $\rightarrow$ policy $\pi$ better or equal to policy $\pi'$

Optimal policy $\pi_*$: best policy that is better or equal to all other policies

Optimal state-value function

$$v_*(s) \doteq \max_\pi v_\pi(s) \qquad \text{for all } s \in \mathcal{S}.$$
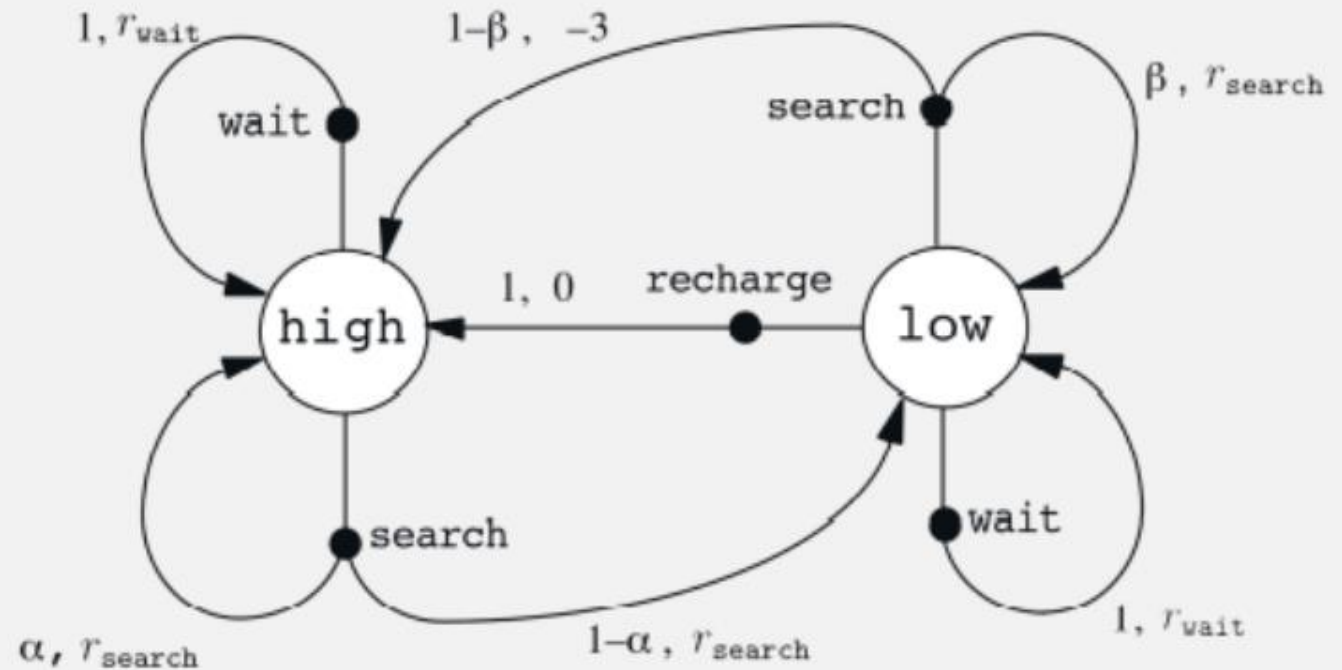
Optimal action-value function

$$q_*(s,a) \doteq \max_\pi q_\pi(s,a) \qquad \text{for all } s \in \mathcal{S} \text{ and } a \in \mathcal{A}(s)$$

$$q_*(s,a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

# Example 3.8: Bellman Optimality Equations for the Recycling robot

| $s$ | $a$ | $s'$ | $p(s'\|s, a)$ | $r(s, a, s')$ |
|------|---------|------|--------------|---------------|
| high | search | high | $\alpha$ | $r_{search}$ |
| high | search | low | $1 - \alpha$ | $r_{search}$ |
| low | search | high | $1 - \beta$ | $-3$ |
| low | search | low | $\beta$ | $r_{search}$ |
| high | wait | high | $1$ | $r_{wait}$ |
| high | wait | low | $0$ | $r_{wait}$ |
| low | wait | high | $0$ | $r_{wait}$ |
| low | wait | low | $1$ | $r_{wait}$ |
| low | recharge | high | $1$ | $0$ |
| low | recharge | low | $0$ | $0$ |

# Example 3.8: cont.

States: {h, l} ~ {high, low}

Actions: {s, w, re} ~ {search, wait, recharge}

Reward: $\{r_s, r_w\}$ ~ $\{r_{search}, r_{wait}\}$

$$v_\pi(s) \doteq \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\left[r + \gamma v_\pi(s')\right], \quad \text{for all } s \in \mathcal{S}$$

$$
\begin{aligned}
v_*(\mathbf{h}) &= \max \left\{
\begin{array}{l}
p(\mathbf{h}|\mathbf{h},\mathbf{s})[r(\mathbf{h},\mathbf{s},\mathbf{h}) + \gamma v_*(\mathbf{h})] + p(\mathbf{l}|\mathbf{h},\mathbf{s})[r(\mathbf{h},\mathbf{s},\mathbf{l}) + \gamma v_*(\mathbf{l})], \\
p(\mathbf{h}|\mathbf{h},\mathbf{w})[r(\mathbf{h},\mathbf{w},\mathbf{h}) + \gamma v_*(\mathbf{h})] + p(\mathbf{l}|\mathbf{h},\mathbf{w})[r(\mathbf{h},\mathbf{w},\mathbf{l}) + \gamma v_*(\mathbf{l})]
\end{array}
\right\} \\
&= \max \left\{
\begin{array}{l}
\alpha[r_{\mathbf{s}} + \gamma v_*(\mathbf{h})] + (1-\alpha)[r_{\mathbf{s}} + \gamma v_*(\mathbf{l})], \\
1[r_{\mathbf{w}} + \gamma v_*(\mathbf{h})] + 0[r_{\mathbf{w}} + \gamma v_*(\mathbf{l})]
\end{array}
\right\} \\
&= \max \left\{
\begin{array}{l}
r_{\mathbf{s}} + \gamma[\alpha v_*(\mathbf{h}) + (1-\alpha)v_*(\mathbf{l})], \\
r_{\mathbf{w}} + \gamma v_*(\mathbf{h})
\end{array}
\right\}.
\end{aligned}
$$

$$
v_*(\mathbf{l}) = \max \left\{
\begin{array}{l}
\beta r_{\mathbf{s}} - 3(1-\beta) + \gamma[(1-\beta)v_*(\mathbf{h}) + \beta v_*(\mathbf{l})] \\
r_{\mathbf{w}} + \gamma v_*(\mathbf{l}), \\
\gamma v_*(\mathbf{h})
\end{array}
\right\}.
$$

For any choice of $r_{\mathbf{s}}$, $r_{\mathbf{w}}$, $\alpha$, $\beta$, and $\gamma$, with $0 \le \gamma < 1$, $0 \le \alpha, \beta \le 1$, there is exactly one pair of numbers, $v_*(\mathbf{h})$ and $v_*(\mathbf{l})$, that simultaneously satisfy these two nonlinear equations. ∎

# Example 3.9: solving the Gridworld



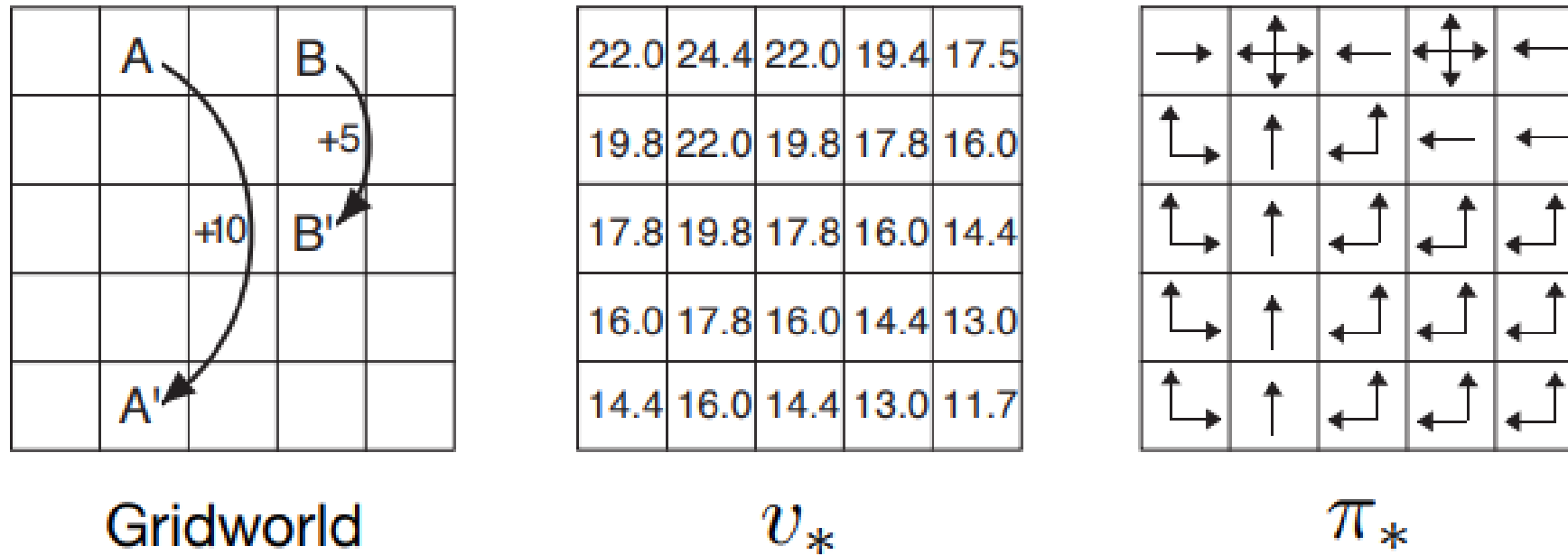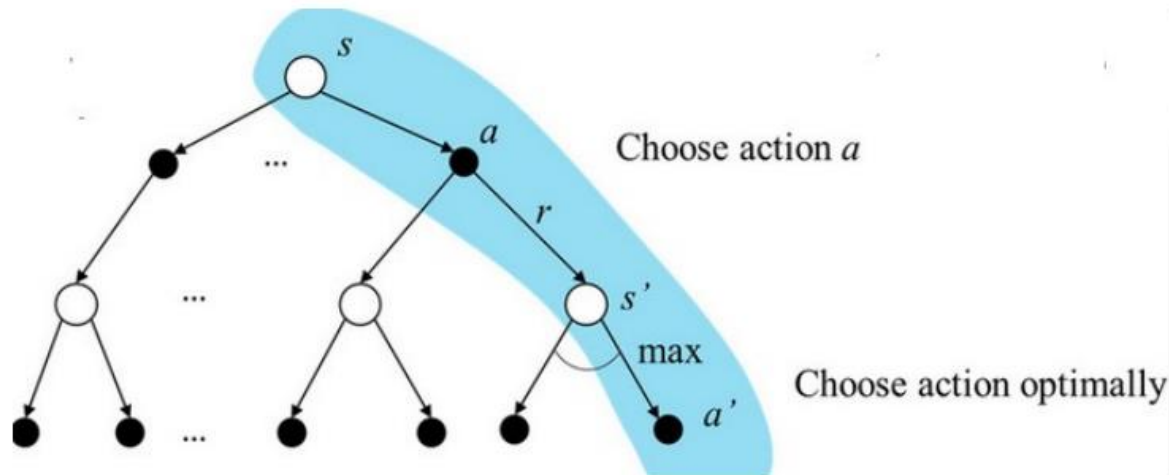Gridworld        $v_*$        $\pi_*$

Figure 3.5: Optimal solutions to the gridworld example.

# Optimal state-value function and Optimal action-value function relationship

$$v_*(s) = \max_a \; q_*(s, a)$$



Choose action $a$

max

Choose action optimally

Nhớ giữ sức khỏe! Take care!

Thank you!
Q&A