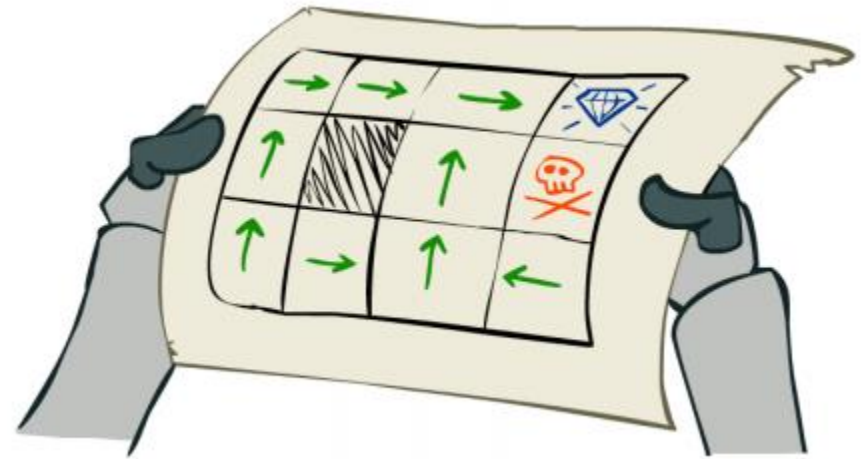
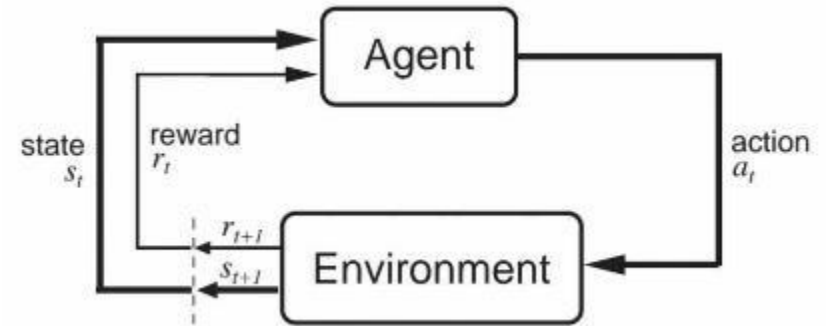


# 가치 반복 알고리즘 Value Iteration



# Infinite-horizon MDP 알고리즘

## • 가치 반복 (value iteration) 알고리즘

$$v^*(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v^*(s') \right\}$$

### • 초기화

- $v_0(s) = 0$  for all  $s$
- $k = 1$

### • 반복 ( $v_{k-1} \rightarrow v_k$ )

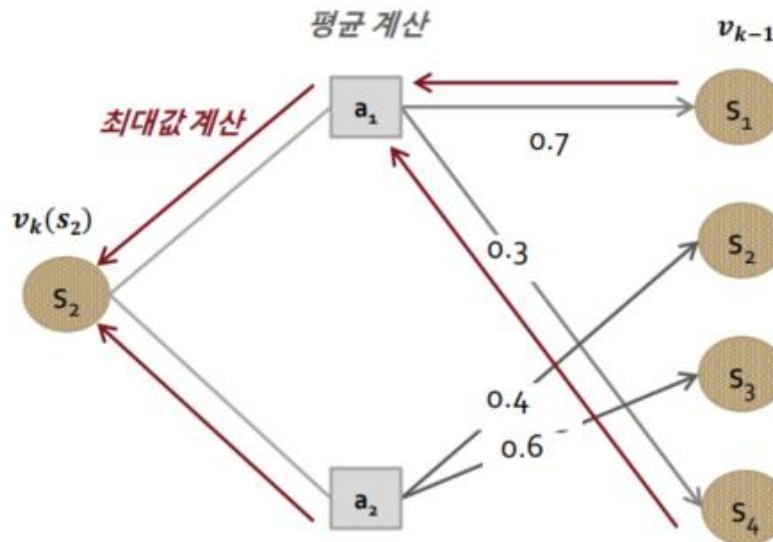
- $v_k(s) = \max_a \{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v_{k-1}(s') \}$
- 만약  $\|v_k - v_{k-1}\| < \epsilon$ , 종료. 아니면  $k \leftarrow k + 1$

# Infinite-horizon MDP 알고리즘

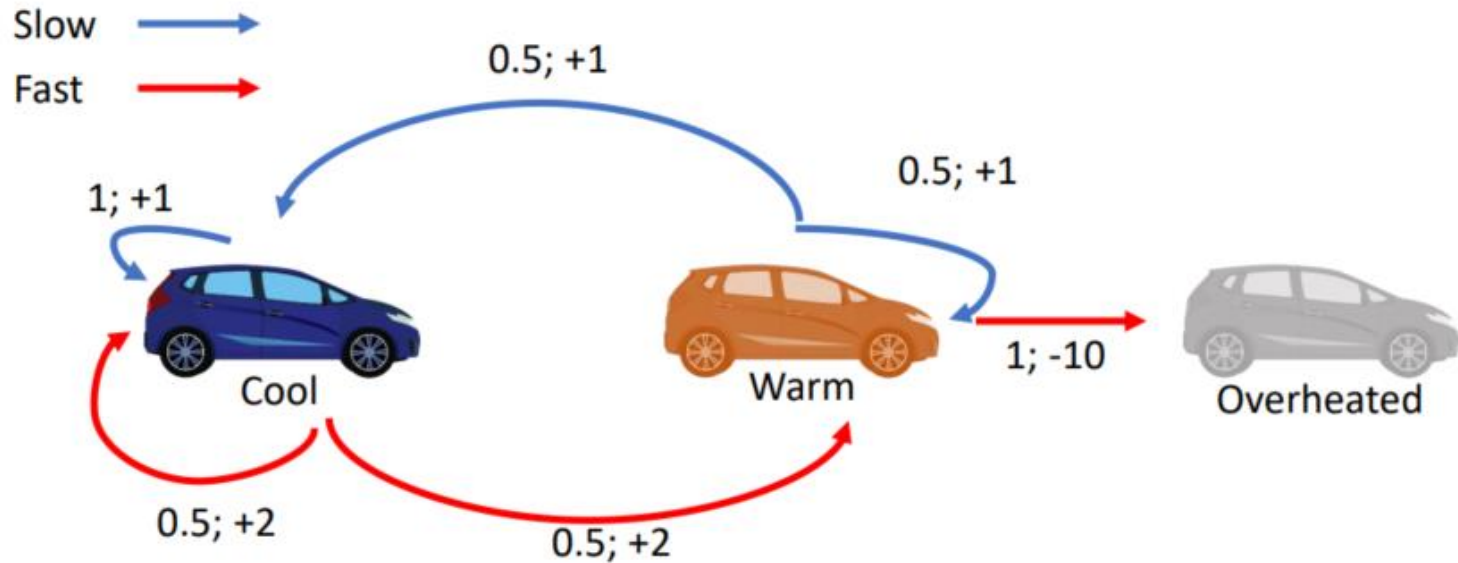
- 가치 반복 (value iteration) 알고리즘

$$v^*(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v^*(s') \right\}$$



$$v_k(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v_{k-1}(s') \right\}$$



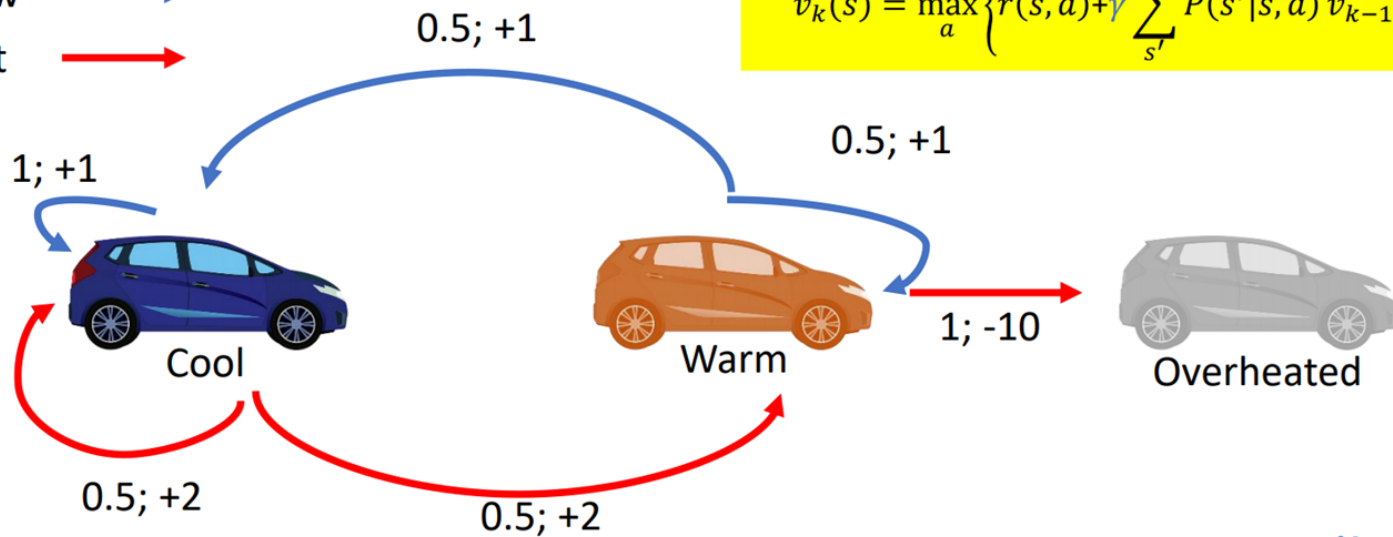
# Infinite-horizon MDP 알고리즘



# Infinite-horizon MDP 알고리즘

Slow   
Fast 

$$v_k(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v_{k-1}(s') \right\}$$



$\gamma = 0.8$

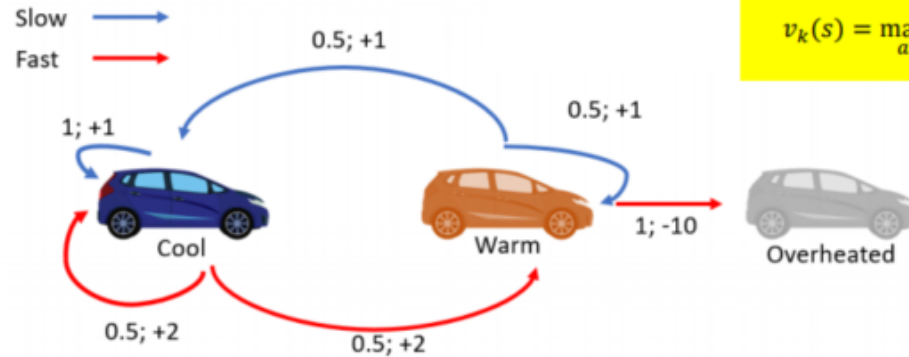
$k$	$v_k(\text{Cool})$	$v_k(\text{Warm})$	$v_k(\text{Overheated})$
0	0	0	0

$$v_1(\text{Cool}) = \max_a \left\{ r(\text{Cool}, a) + \gamma \sum_{s'} P(s'|\text{Cool}, a) v_0(s') \right\} = \max\{r(\text{Cool}, \text{Slow}), r(\text{Cool}, \text{Fast})\} = \max\{+1, +2\} = +2$$

$$v_1(\text{Warm}) = \max_a \left\{ r(\text{Warm}, a) + \gamma \sum_{s'} P(s'|\text{Warm}, a) v_0(s') \right\} = \max\{r(\text{Warm}, \text{Slow}), r(\text{Warm}, \text{Fast})\} = \max\{+1, -10\} = +1$$

$$v_1(\text{Overheated}) = 0$$

# Infinite-horizon MDP 알고리즘



$$v_k(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v_{k-1}(s') \right\}$$

$\gamma = 0.8$

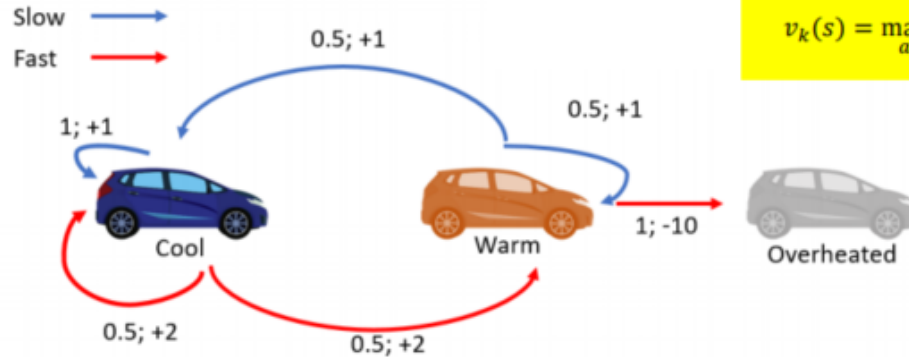
$k$	$v_k(\text{Cool})$	$v_k(\text{Warm})$	$v_k(\text{Overheated})$
0	0	0	0
1	2	1	0

$$v_2(\text{Cool}) = \max_a \left\{ r(\text{Cool}, a) + \gamma \sum_{s'} P(s'|\text{Cool}, a) v_1(s') \right\} = \max \left\{ \begin{array}{l} r(\text{Cool}, \text{Slow}) + 0.8v_1(\text{Cool}) \\ r(\text{Cool}, \text{Fast}) + 0.8(0.5v_1(\text{Cool}) + 0.5v_1(\text{Warm})) \end{array} \right\} = \max \left\{ \begin{array}{l} 1 + 0.8 \times 2 \\ 2 + 0.8(0.5 \times 2 + 0.5 \times 1) \end{array} \right\} = 3.2$$

$$v_2(\text{Warm}) = \max_a \left\{ r(\text{Warm}, a) + \gamma \sum_{s'} P(s'|\text{Warm}, a) v_1(s') \right\} = \max \left\{ \begin{array}{l} r(\text{Warm}, \text{Slow}) + 0.8(0.5v_1(\text{Cool}) + 0.5v_1(\text{Warm})) \\ r(\text{Warm}, \text{Fast}) + 0.8v_1(\text{Overheated}) \end{array} \right\} \\ = \max \left\{ \begin{array}{l} 1 + 0.8(0.5 \times 2 + 0.5 \times 1) \\ -10 + 0.8 \times 0 \end{array} \right\} = 2.2$$

$$v_2(\text{Overheated}) = 0$$

# Infinite-horizon MDP 알고리즘



$$v_k(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v_{k-1}(s') \right\}$$

$\gamma = 0.8$

$k$	$v_k(\text{Cool})$	$v_k(\text{Warm})$	$v_k(\text{Overheated})$
0	0	0	0
1	2	1	0
2	3.2	2.2	0

$$v_3(\text{Cool}) = \max_a \left\{ r(\text{Cool}, a) + \gamma \sum_{s'} P(s'|\text{Cool}, a) v_2(s') \right\} = \max \left\{ \begin{array}{l} r(\text{Cool}, \text{Slow}) + 0.8v_2(\text{Cool}) \\ r(\text{Cool}, \text{Fast}) + 0.8(0.5v_2(\text{Cool}) + 0.5v_2(\text{Warm})) \end{array} \right\}$$

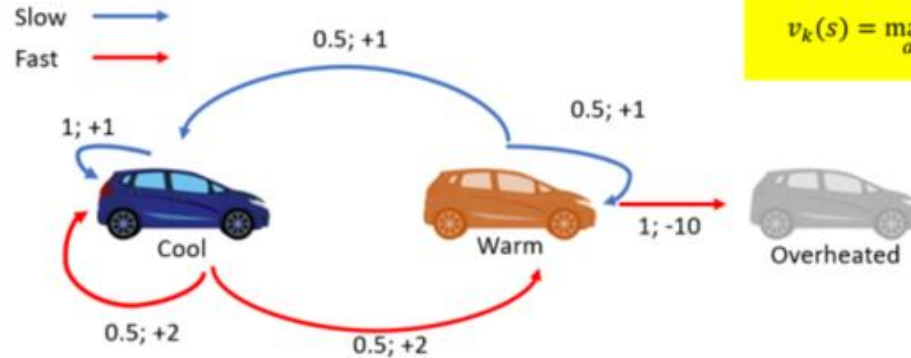
$$= \max \left\{ \begin{array}{l} 1 + 0.8 \times 3.2 \\ 2 + 0.8(0.5 \times 3.2 + 0.5 \times 2.2) \end{array} \right\} = 4.16$$

$$v_3(\text{Warm}) = \max_a \left\{ r(\text{Warm}, a) + \gamma \sum_{s'} P(s'|\text{Warm}, a) v_2(s') \right\} = \max \left\{ \begin{array}{l} r(\text{Warm}, \text{Slow}) + 0.8(0.5v_2(\text{Cool}) + 0.5v_2(\text{Warm})) \\ r(\text{Warm}, \text{Fast}) + 0.8v_2(\text{Overheated}) \end{array} \right\}$$

$$= \max \left\{ \begin{array}{l} 1 + 0.8(0.5 \times 3.2 + 0.5 \times 2.2) \\ -10 + 0.8 \times 0 \end{array} \right\} = 3.16$$

$$v_3(\text{Overheated}) = 0$$

# Infinite-horizon MDP 알고리즘



$$v_k(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} P(s'|s, a) v_{k-1}(s') \right\}$$

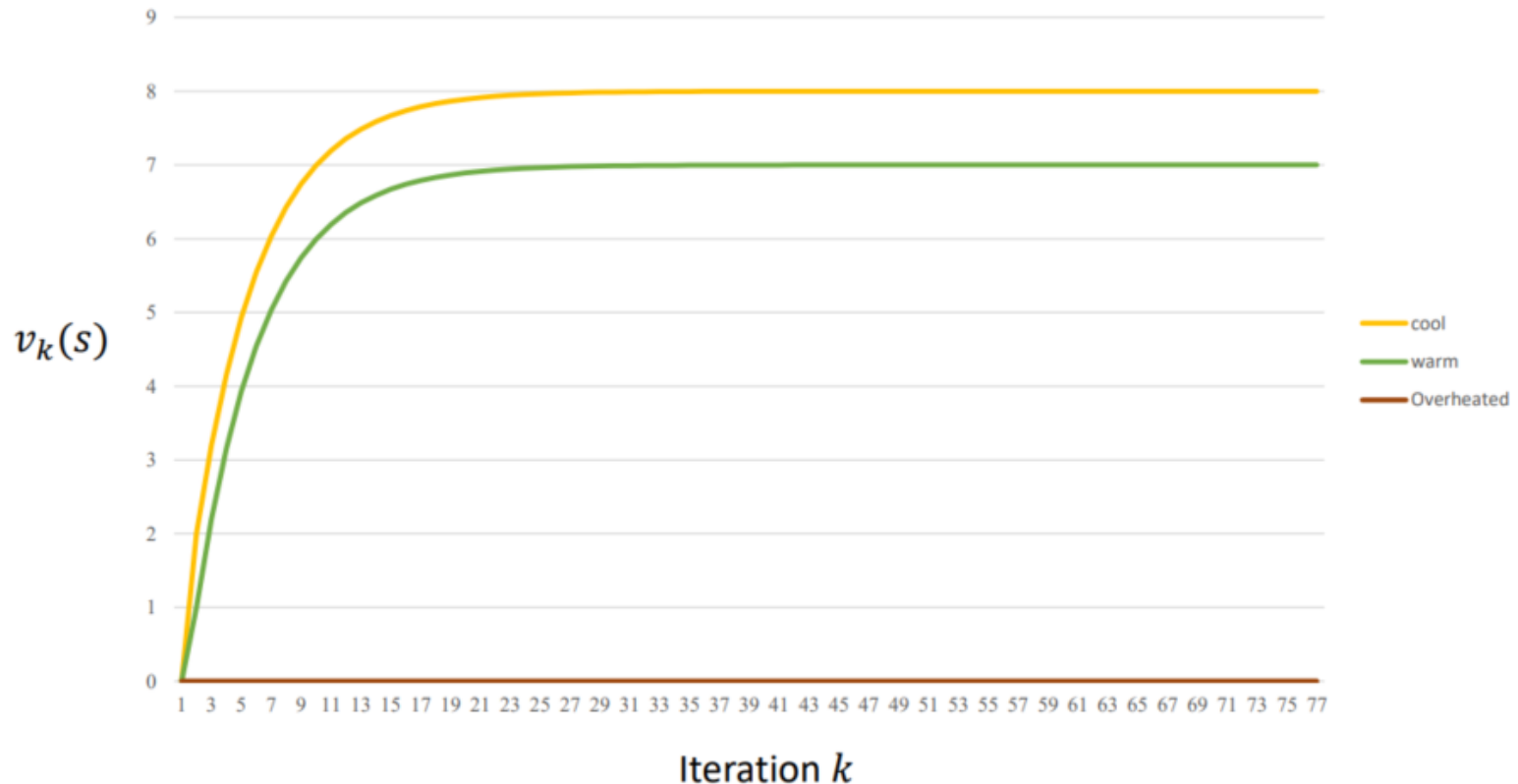
$\gamma = 0.8$

$k$	$v_k(\text{Cool})$	$v_k(\text{Warm})$	$v_k(\text{Overheated})$
0	0	0	0
1	2	1	0
2	3.2	2.2	0
3	4.16	3.16	0

⋮



# Infinite-horizon MDP 알고리즘



# Infinite-horizon MDP 알고리즘

## • 최적 정책은?

$$\bullet \pi^*(s) = \operatorname{argmax}_a \{r(s, a) + \gamma \sum_{s'} P(s'|s, a) v^*(s')\}$$

$$\bullet \pi^*(Cool) = \operatorname{argmax} \left\{ \begin{array}{l} 1 + 0.8 \times v^*(Cool) \\ 2 + 0.8(0.5v^*(Cool) + 0.5v^*(Warm)) \end{array} \right\}$$

$$= \operatorname{argmax} \left\{ \begin{array}{l} 1 + 0.8 \times 8 \\ 2 + 0.8(0.5 \times 8 + 0.5 \times 7) \end{array} \right\} = Fast$$

$$\bullet \pi^*(Warm) = \operatorname{argmax} \left\{ \begin{array}{l} 1 + 0.8(0.5v^*(Cool) + 0.5v^*(Warm)) \\ -10 + 0.8v^*(Overheated) \end{array} \right\}$$

$$= \operatorname{argmax} \left\{ \begin{array}{l} 1 + 0.8(0.5 \times 8 + 0.5 \times 7) \\ -10 + 0.8 \times 0 \end{array} \right\} = Slow$$

