



# Markov Decision Process

(MDP)

Edgar Granados

# R

# Sequential Decision Making

## Definition

- Episodic Decisions...
- Continuous decision making?

# R

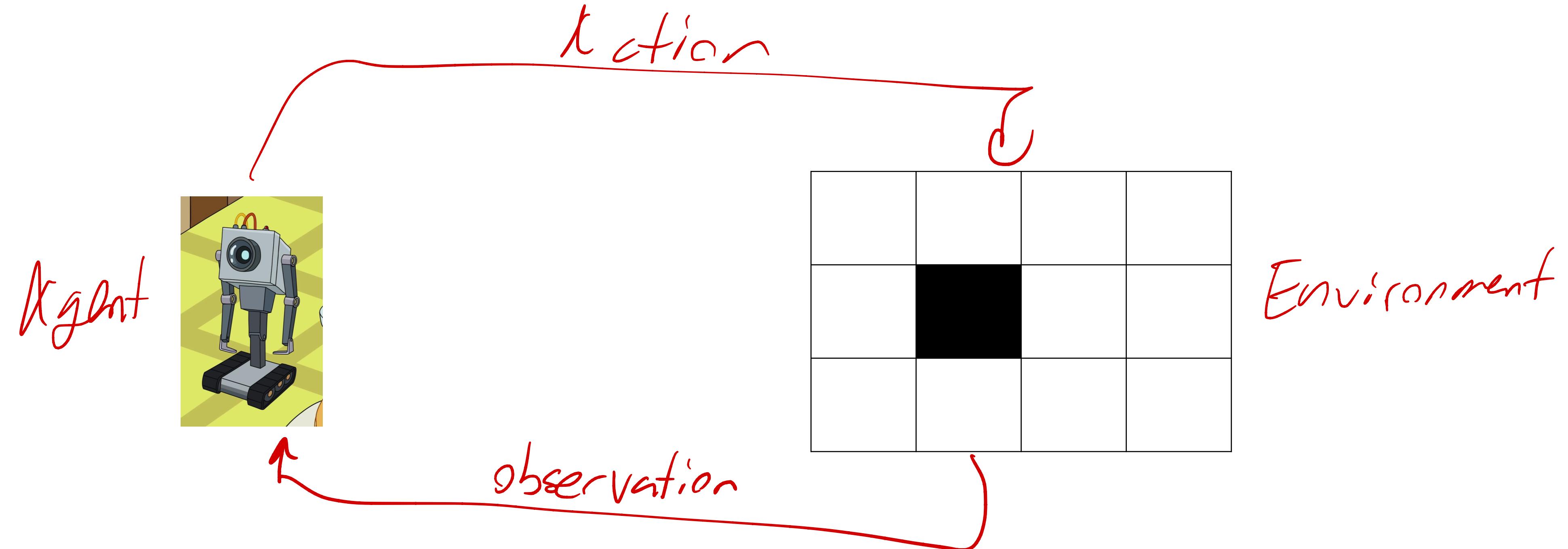
# Sequential Decision Making

## Definition

- Episodic Decisions...
  - Continuous decision making?
  - **Take actions that affects later decisions**
- 

# R

# Agent-Environment Interaction



- Probabilistic actions
- Deterministic observations

- Sequence that yields the best payout - Reward

# R

# Agent-Environment Interaction

## Example

- Blackjack

10    9  
18

13

A - K  
A



# R

## Markov Decision Process

### Definition

- Tuple  $(\underline{S}, \underline{A}, \underline{T}, \underline{R})$ :
  - $S$  - States
  - $A$  - Actions
  - $T$  - Transition Matrix
  - $R$  - Reward function

Policy  $\pi^*: S \rightarrow A$

**R**

# MDPs

**Markov?**

- Markov Assumption!
- Returns a **Policy**



# R

## MDP

### States

- States - Representation of all ~~reveal~~<sup>relevant</sup> information for predicting future states
  - Describes the configuration of the system at a given moment



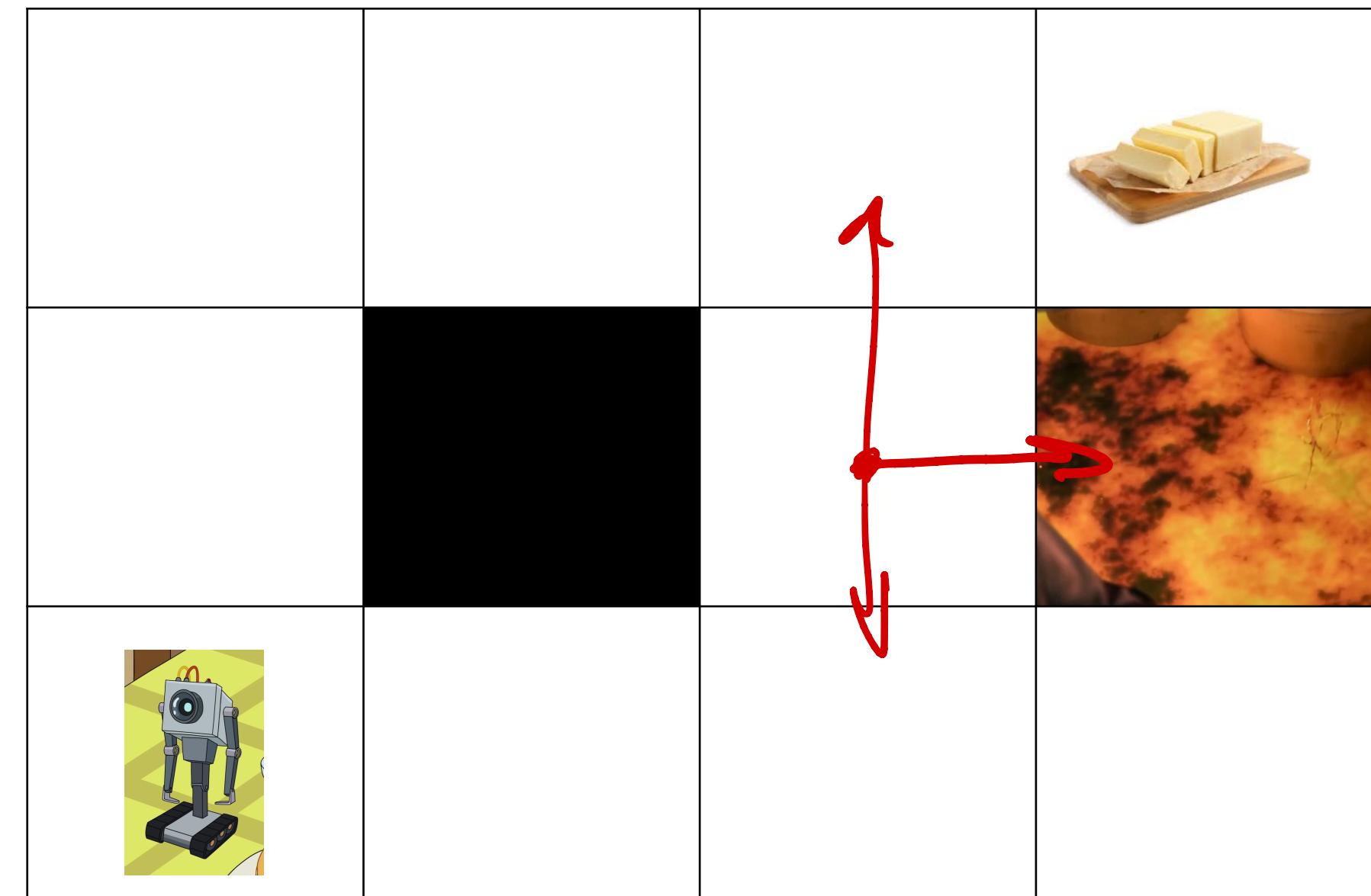
# R

## MDP

### Actions

- When executed, actions modify states
- Determine the best sequence of actions

$$A = \{L, R, U, D\}$$



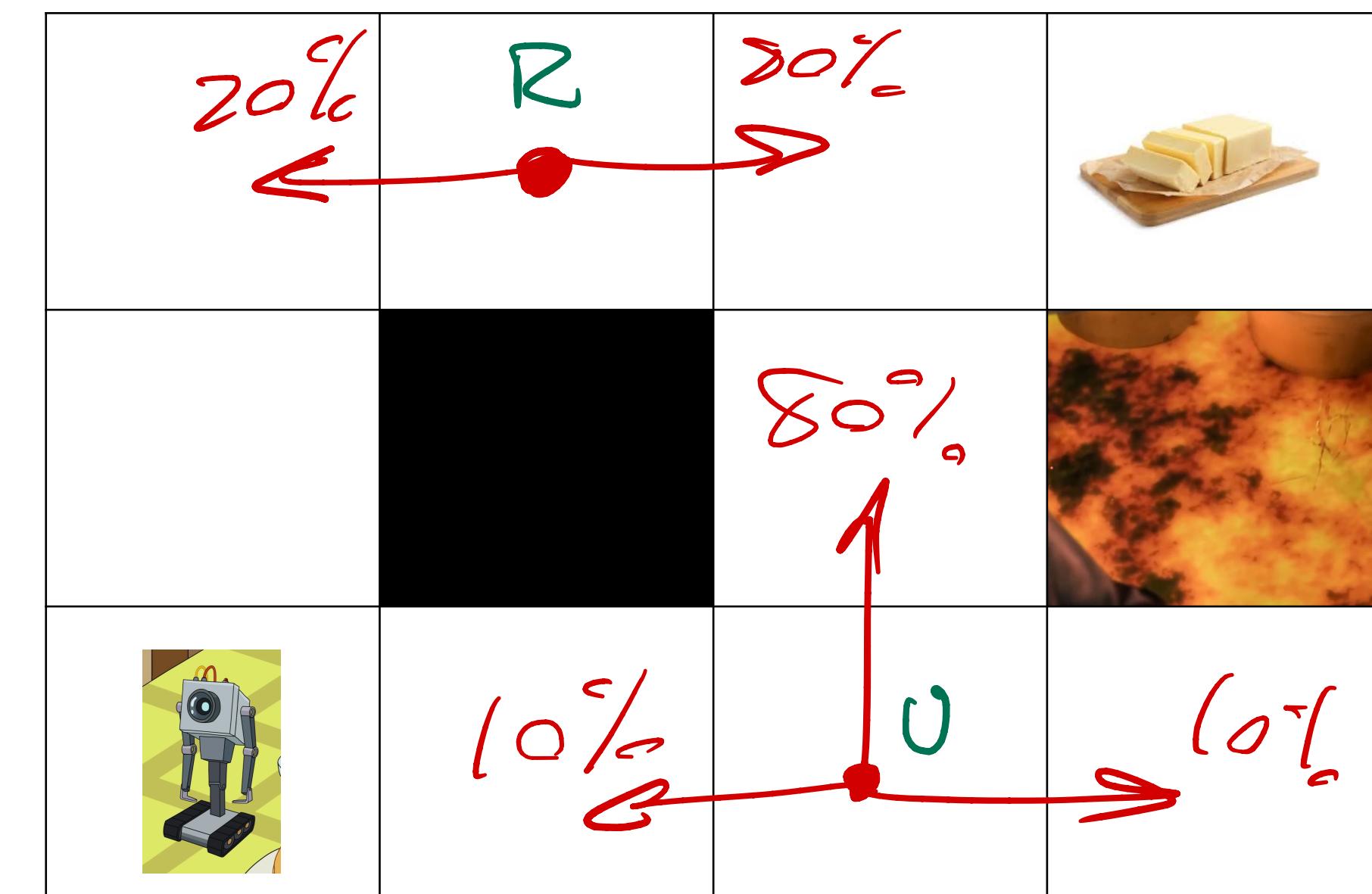
# R

## MDP

### Transition Function

- $T(\underline{s}_t, a_t, \underline{s}_{t+1}) = P(\underline{s}_{t+1} | s_t, a_t)$

$80\%$   $\rightarrow$  correct state | action  
 $20\%$   $\rightarrow$  incorrect state(s)



# R

## MDP

### Reward Function

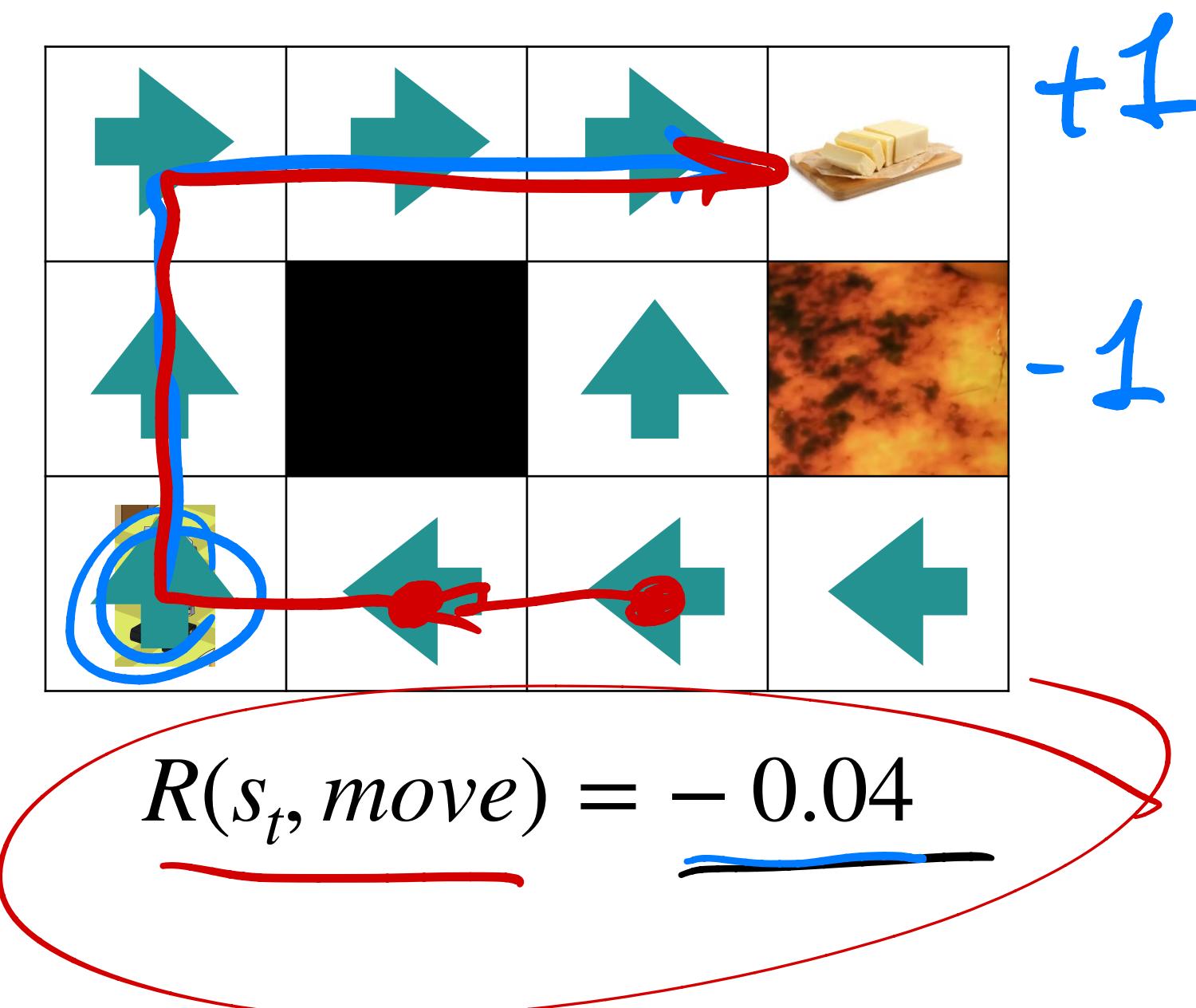
- $R(s_t, a_t)$
- Direct agent towards desirable states & away from unwanted ones
- Maximize the cumulative reward
- MDP: Model complex concurrent tasks by simply assigning rewards to the states

			+1 
			-1 

# R

# Rewards

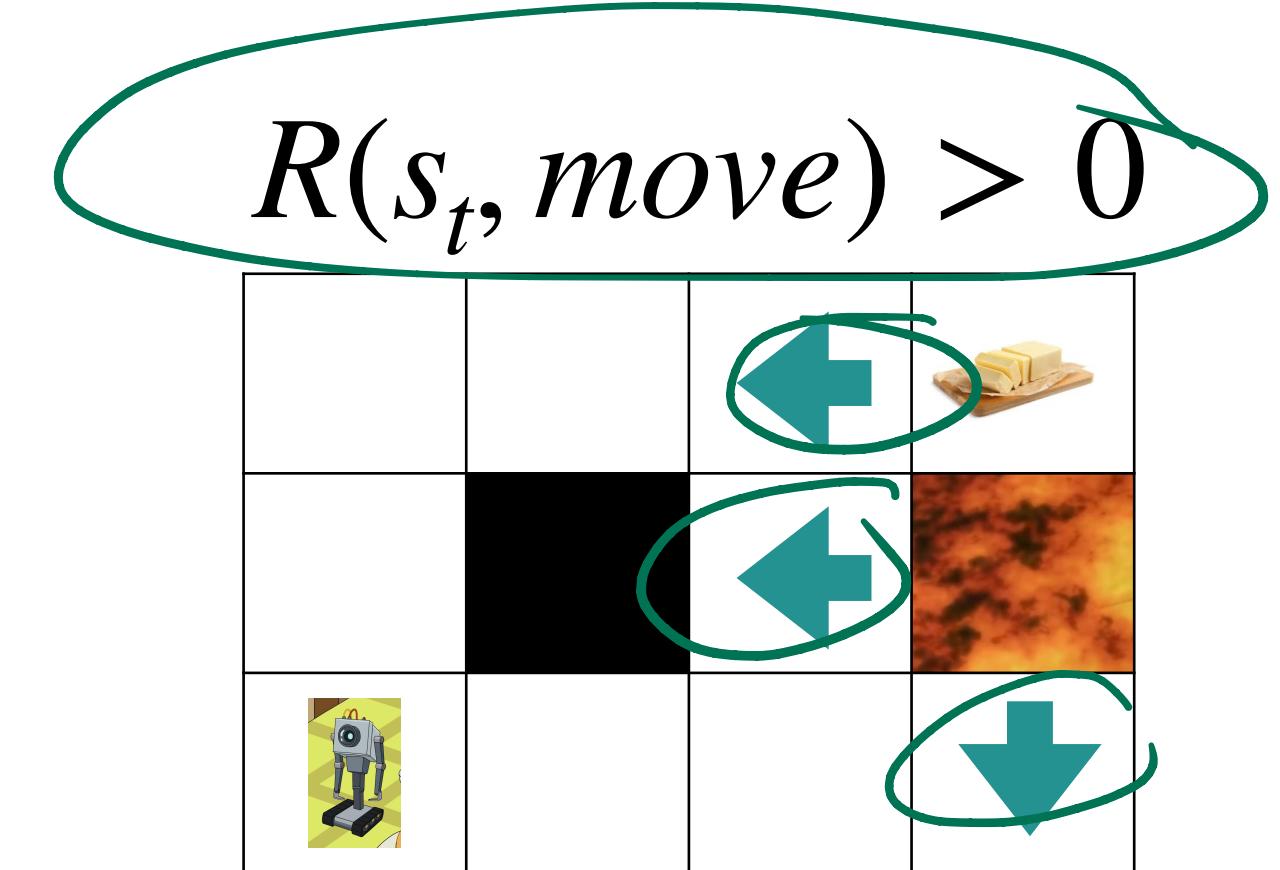
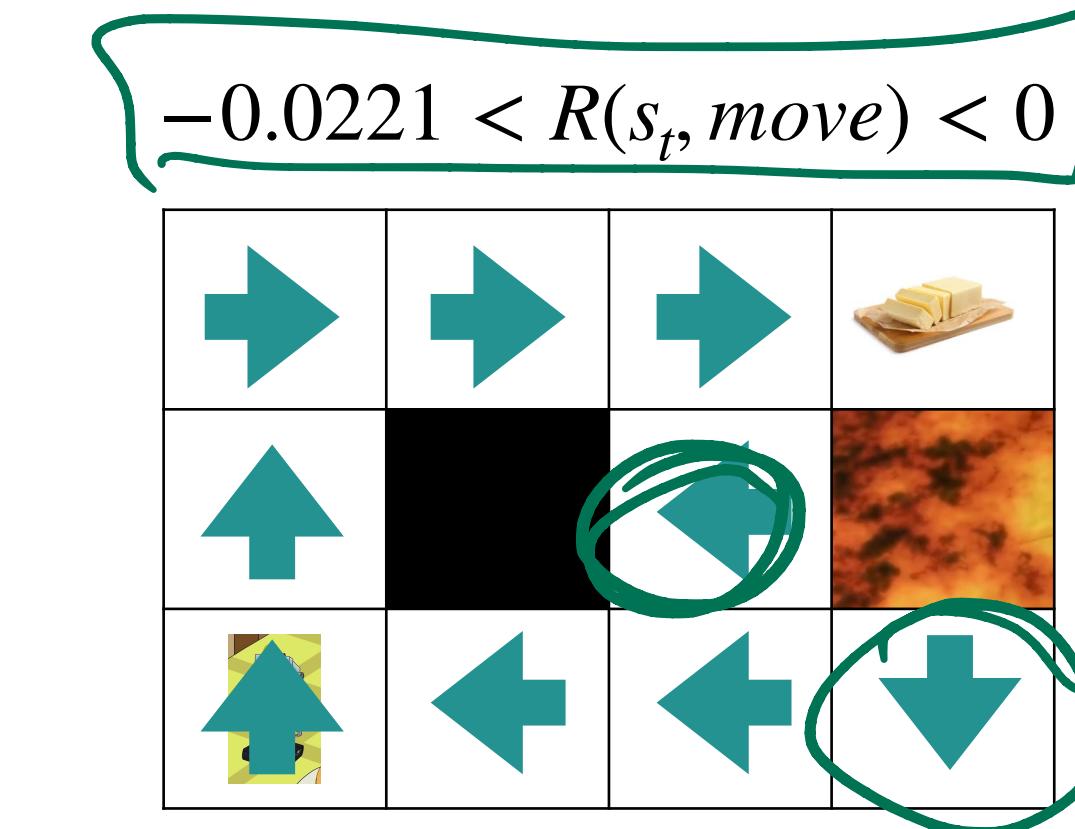
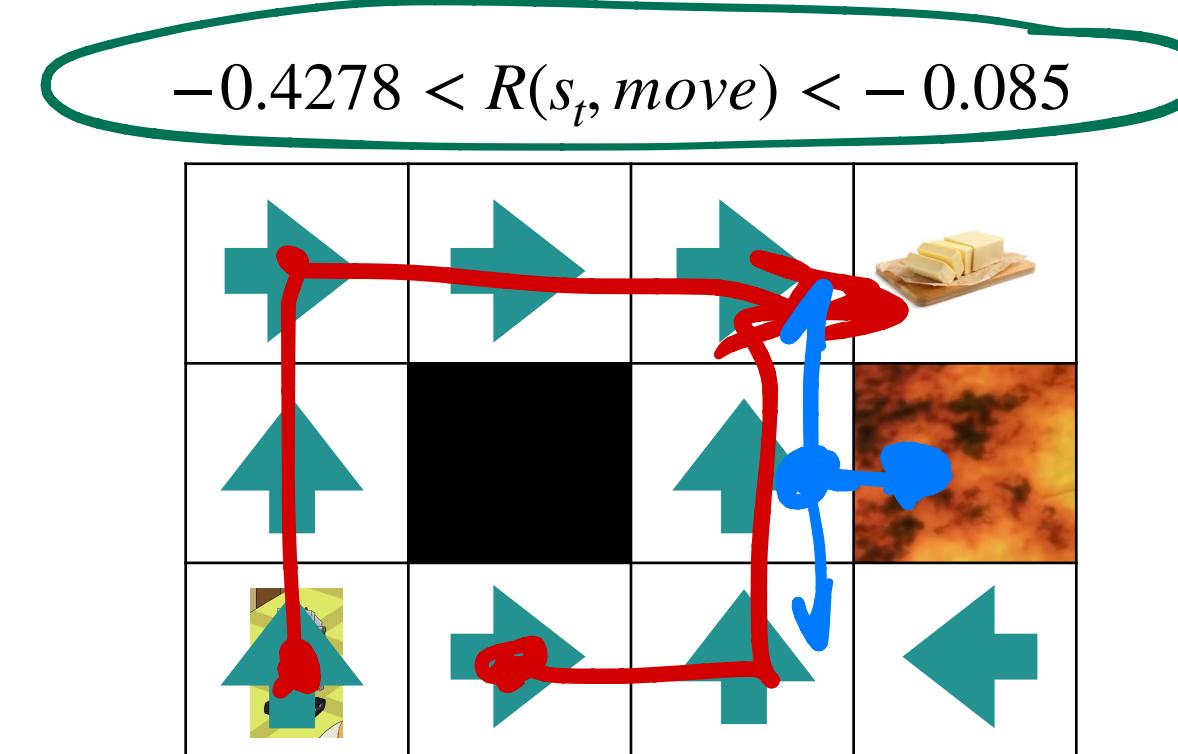
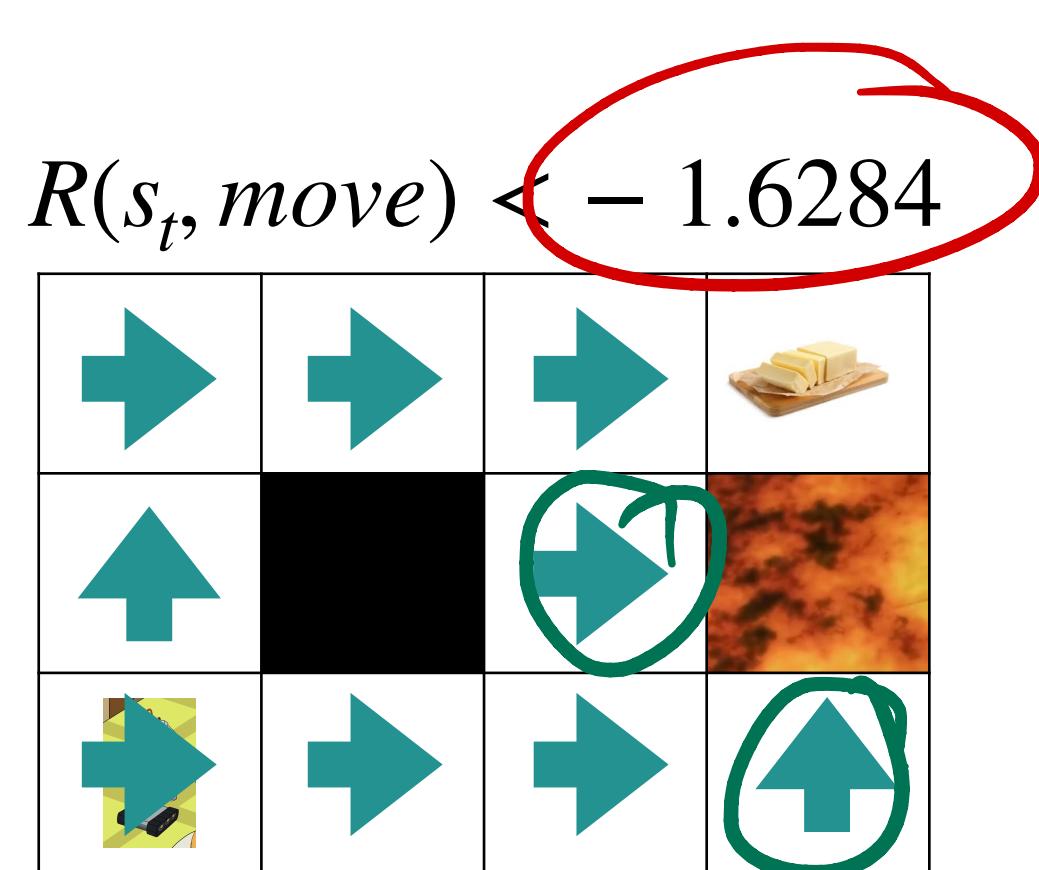
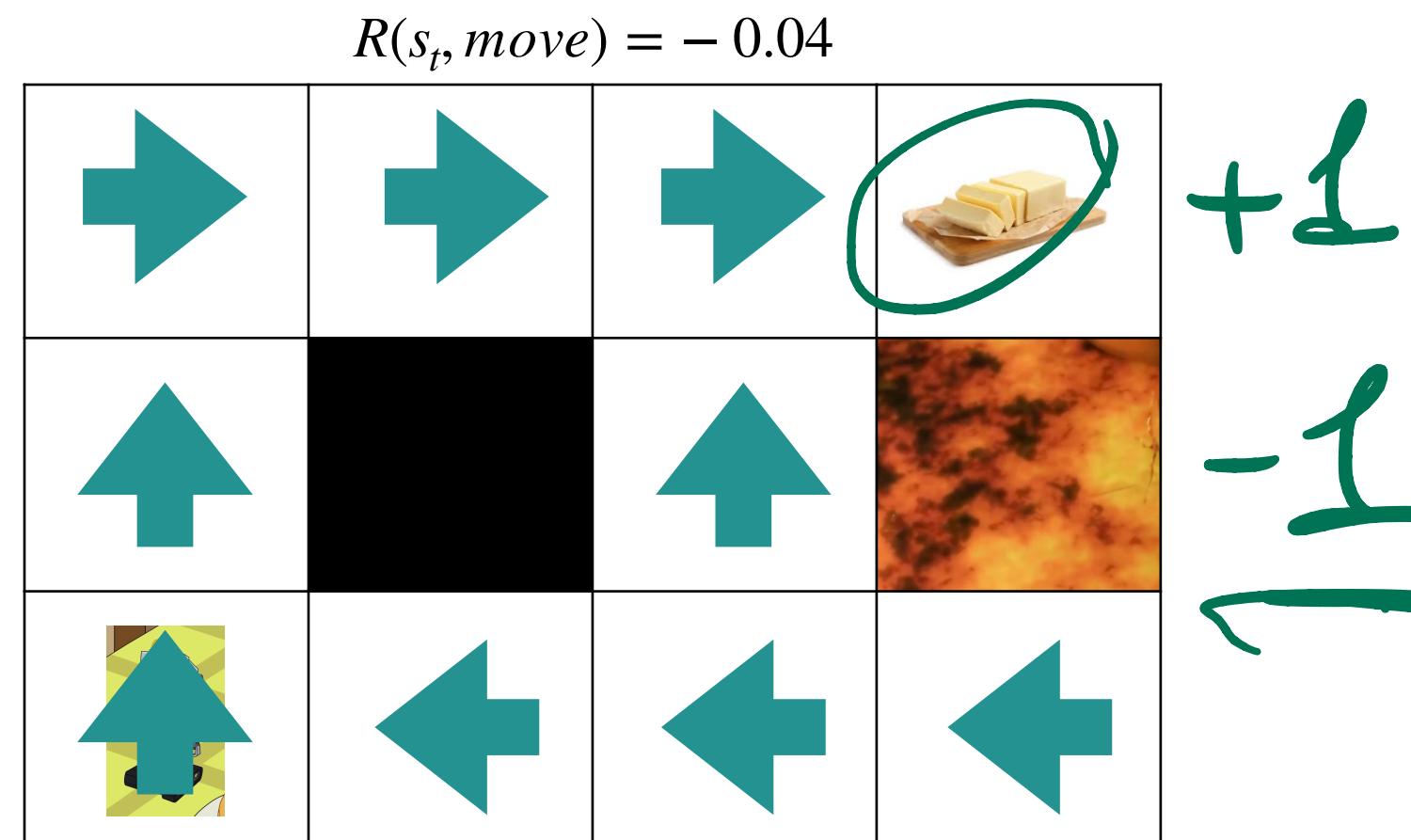
*Appropriate rewards?*



# R

# Rewards

*Appropriate rewards?*



# R

## MDP

- Overall utility:

- $V[s_0, s_1, \dots] = \overbrace{R(s_0) + R(s_1) + \dots}^{\text{red underline}}$

# R

## MDP

- Overall utility using discounts:  
$$V[s_0, s_1, \dots] = \underbrace{R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \gamma^3 R(s_3) + \dots}_{\text{0} < \gamma < 1}$$

$$\frac{R_{max}}{1-\gamma} \rightarrow \text{Max Possible total Reward}$$

# R

## Bellman's Equation

- $\pi : S \rightarrow A:$

$$\begin{aligned} V^\pi(s) &= \sum_{t=0}^{\infty} \gamma^t E[R(s_t) | s_0 = s] = R(s) + \sum_{t=1}^{\infty} \gamma^t E[R(s_t) | s_0 = s] \\ &= R(s) + \gamma \sum_{t=0}^{\infty} \gamma^t E[R(s) | s_0 \in T(s, \pi(s))] \\ &= R(s) + \gamma \sum_{s' \in T(s, \pi(s))} P(s'|s, \pi(s)) \sum_{t=0}^{\infty} \gamma^t E[R(s_t) | s_0 = s] \end{aligned}$$

**R**

# Bellman's Equation

$$V^\pi(s) = \underbrace{R(s, \pi(s))}_{\alpha} + \gamma \sum_{s' \in S} T(s, \pi(s), s') V^\pi(s')$$

# R

## Bellman's Equation

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} T(s, \pi(s), s') V^\pi(s')$$

- Dynamic programming:
  - Solve complex problems by breaking them down into simpler ones

**R**

# Optimal Policy

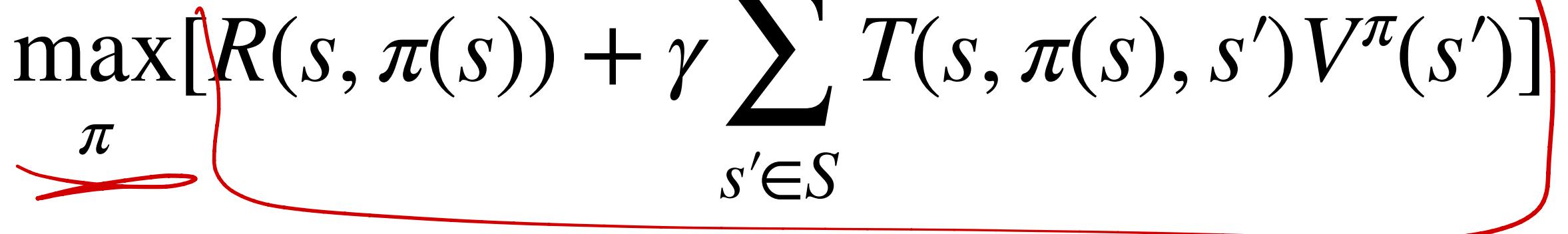
 $\pi^*$ 

- An optimal policy satisfies:
  - $\forall s \in S : \pi^* \in \arg \max V^\pi(s)$

# R

# Optimal Policy

- An optimal policy satisfies:
  - $\forall s \in S : \underline{\pi^*} \in \arg \max V^\pi(s)$
- The necessary and sufficient condition:

$$V^\pi(s) = \max_{\pi} [R(s, \pi(s)) + \gamma \sum_{s' \in S} T(s, \pi(s), s') V^\pi(s')]$$


# R

$$V_{k+1}(s) = R(s) + \gamma \max_{a \in A(S)} \left[ \sum_{s' \in T(s,a)} P(s'|s,a) V_k(s') \right]$$

1.  $V_0 \equiv 0, k = 0$

2. Compute  $v_{k+1}, \forall s \in S$

3. (Optional) if  $|V_{k+1}(s) - V_k(s)| > \delta$ :

1.  $|V_{k+1}(s) - V_k(s)| = \delta$

4. Repeat until  $|V_{k+1}(s) - V_k(s)| < \epsilon \forall s \in S$

$$\epsilon > 0$$

# Solving MDPs

## Value Iteration

$$N = \left\lceil \frac{\log \left( \frac{2R_{\max}}{\epsilon(1-\gamma)} \right)}{\log(1-\gamma)} \right\rceil$$

$\uparrow \gamma \Rightarrow \text{Slow}$

# R

$$V_{k+1}(s) = R(s) + \gamma \max_{a \in A(S)} \left[ \sum_{s' \in T(s,a)} P(s'|s,a) V_k(s') \right]$$

(Handwritten annotations: A red bracket highlights the term  $\sum_{s' \in T(s,a)} P(s'|s,a) V_k(s')$ . A green oval encloses the term  $P(s'|s,a) V_k(s')$ .

# Solving MDPs

## Value Iteration

$$\gamma = 1$$

1.  $V_0([4,3]) = 1, V_0([4,2]) = -1$

$$V_0(s') = \emptyset$$

2. Start at [3,3]  $\sum_{s' \in T(s,a)} P(s'|s,a) V_0(s')$

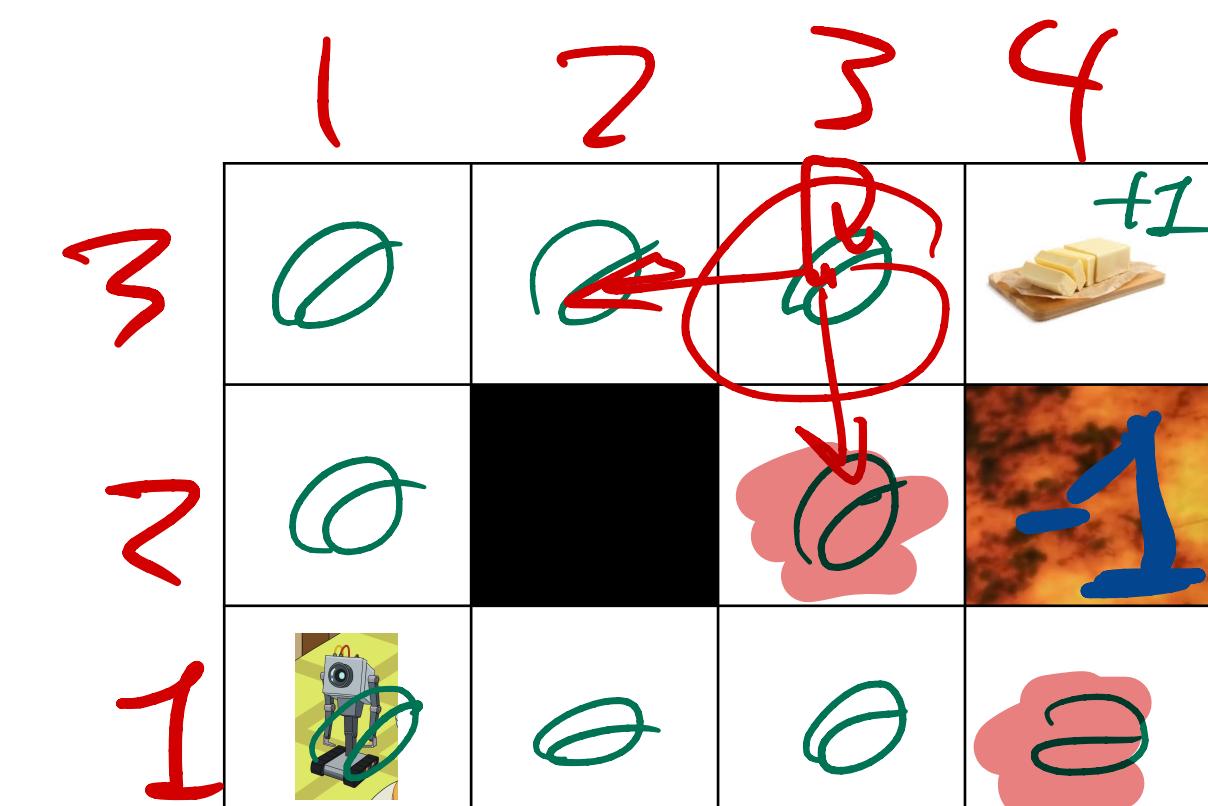
1. Up:  $.8(0) + .1(0) + .1(1) = 1$

2. Right:  $.8(1) + .1(0) + .1(0) = 0.8$

3. Down:  $.8(0) + .1(0) + .1(1) = 1$

4. Left:  $.8(0) + .1(0) + .1(0) = 0$

$\max$



$$V_1([3,3]) = \underline{-0.04} + .8 = .76$$

Right

# R

$$V_{k+1}(s) = R(s) + \gamma \max_{a \in A(S)} \left[ \sum_{s' \in T(s,a)} P(s'|s,a) V_k(s') \right]$$

# Solving MDPs

## Value Iteration

1.  $V_0([4,3]) = 1, V_0([4,2]) = -1$
2. Start at [3,3]  $\sum_{s' \in T(s,a)} P(s'|s,a) V_0(s')$ 
  1. Best action is going right

3. All other  $s \in S, V_1(s) = -0.04$

4. At [3,3], compute  $V_2$

1. Up  $.8(.76) + .1(1) + .1(-.04) = .668$

2. Right  $.8(1) + .1(-.04) + .1(.76) = .872$

3. Down  $.8(-.04) + .1(1) + .1(-.04) = -.064$

4. Left  $.8(-.04) + .1(-.04) + .1(.76) = .039$

$$V_2([3,3]) = -0.04 + .872 = .832$$

$-.04$	$-.04$	$.76$	
$-.04$		$-.04$	
	$-.04$	$-.04$	$-.04$

Go Right

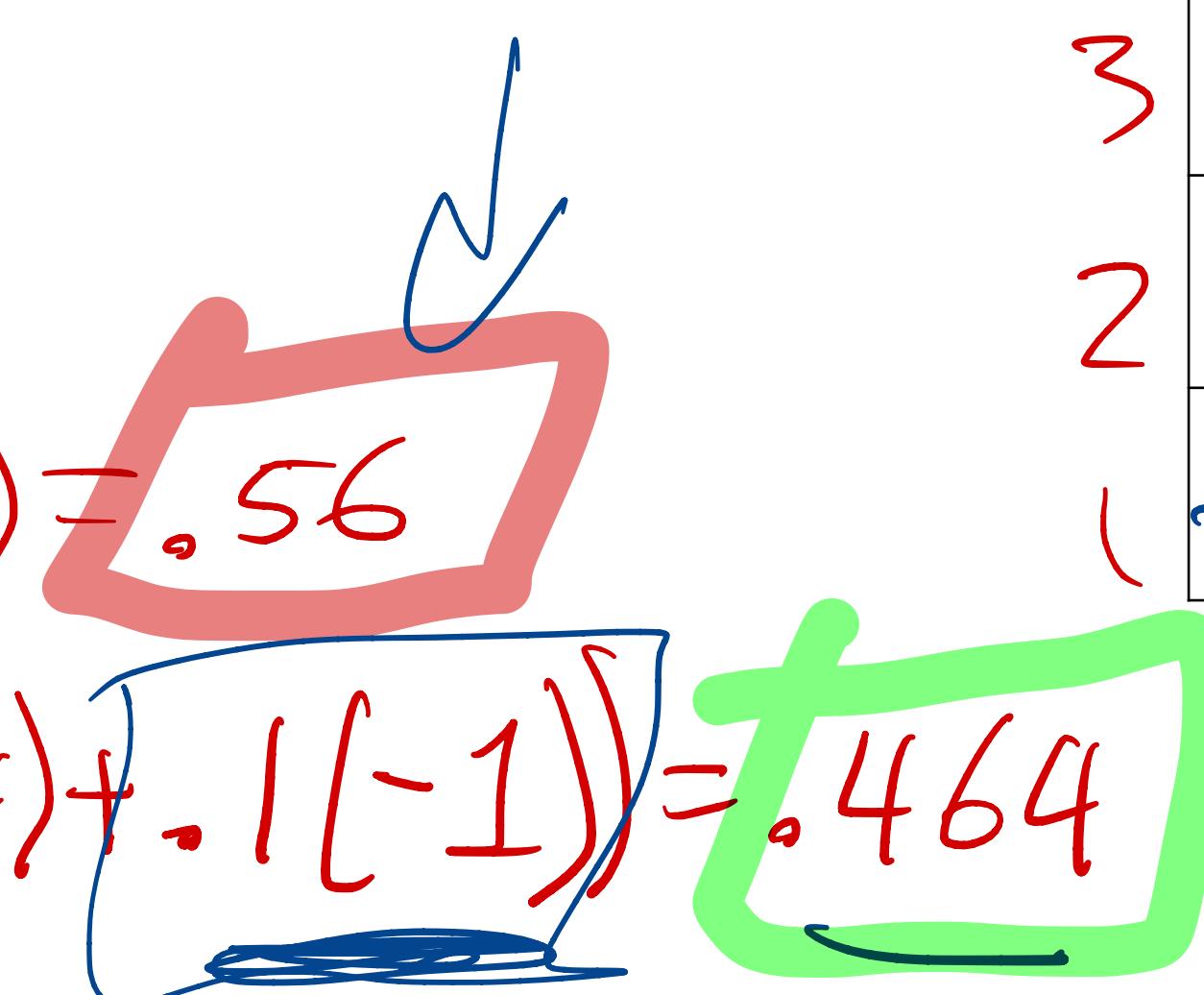
# R

$$V_{k+1}(s) = R(s) + \gamma \max_{a \in A(S)} \left[ \sum_{s' \in T(s,a)} P(s'|s,a) V_k(s') \right]$$

# Solving MDPs

## Value Iteration

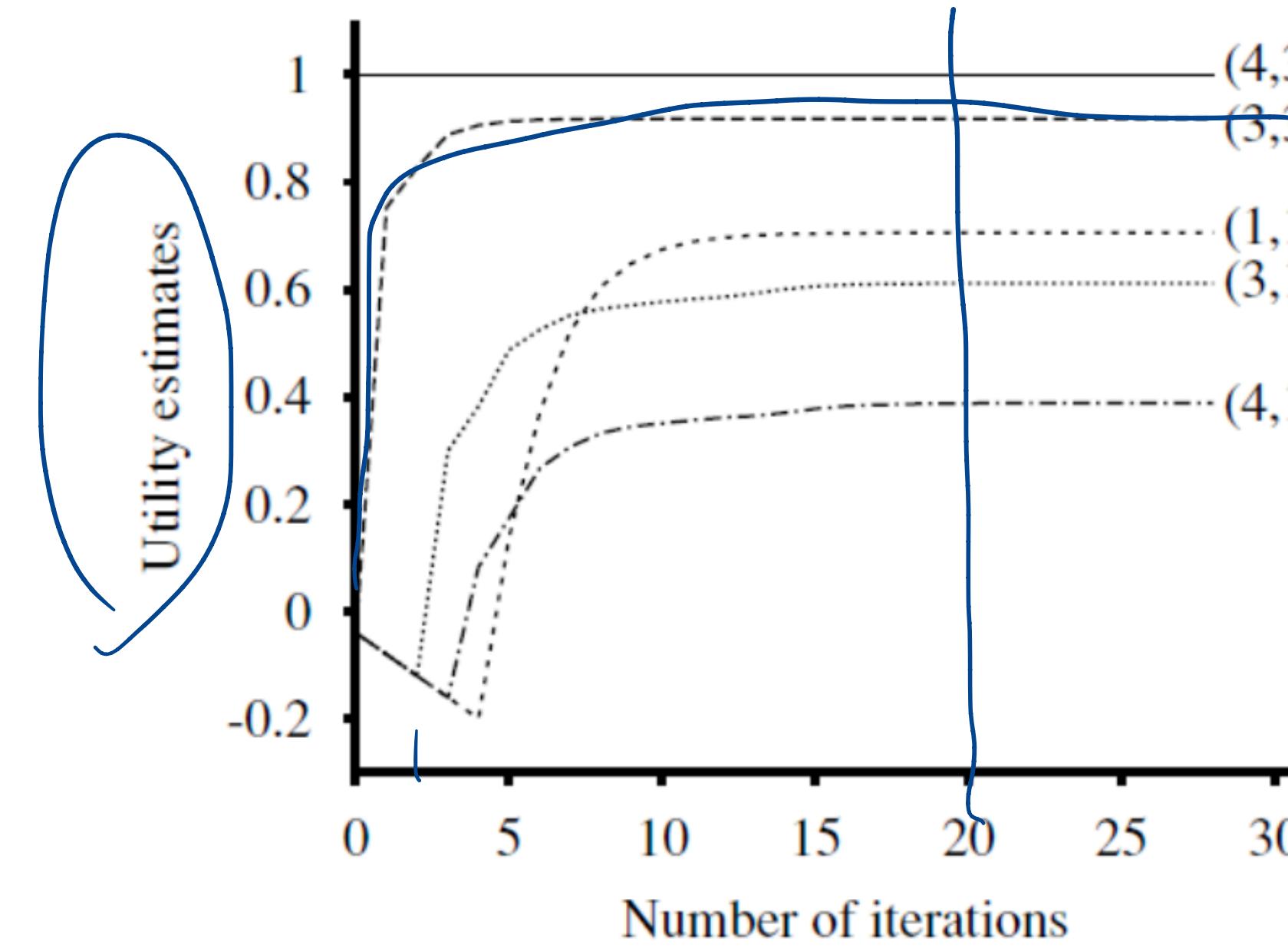
1.  $V_0([4,3]) = 1, V_0([4,2]) = -1$
2. Start at [3,3]  $\sum_{s' \in T(s,a)} P(s'|s,a) V_0(s')$ 
  1. Best action is going **right**
3. All other  $s \in S, V_1(s) = -0.04$
4. At [3,3], compute  $V_2$
5.  $V_2([2,3]) = -.04 + (.8(-.76) + .2(-.04)) = -.56$
6.  $V_2([3,2]) = -.04 + (.8(-.76) + .1(-.04) + .1(-1)) = -.464$



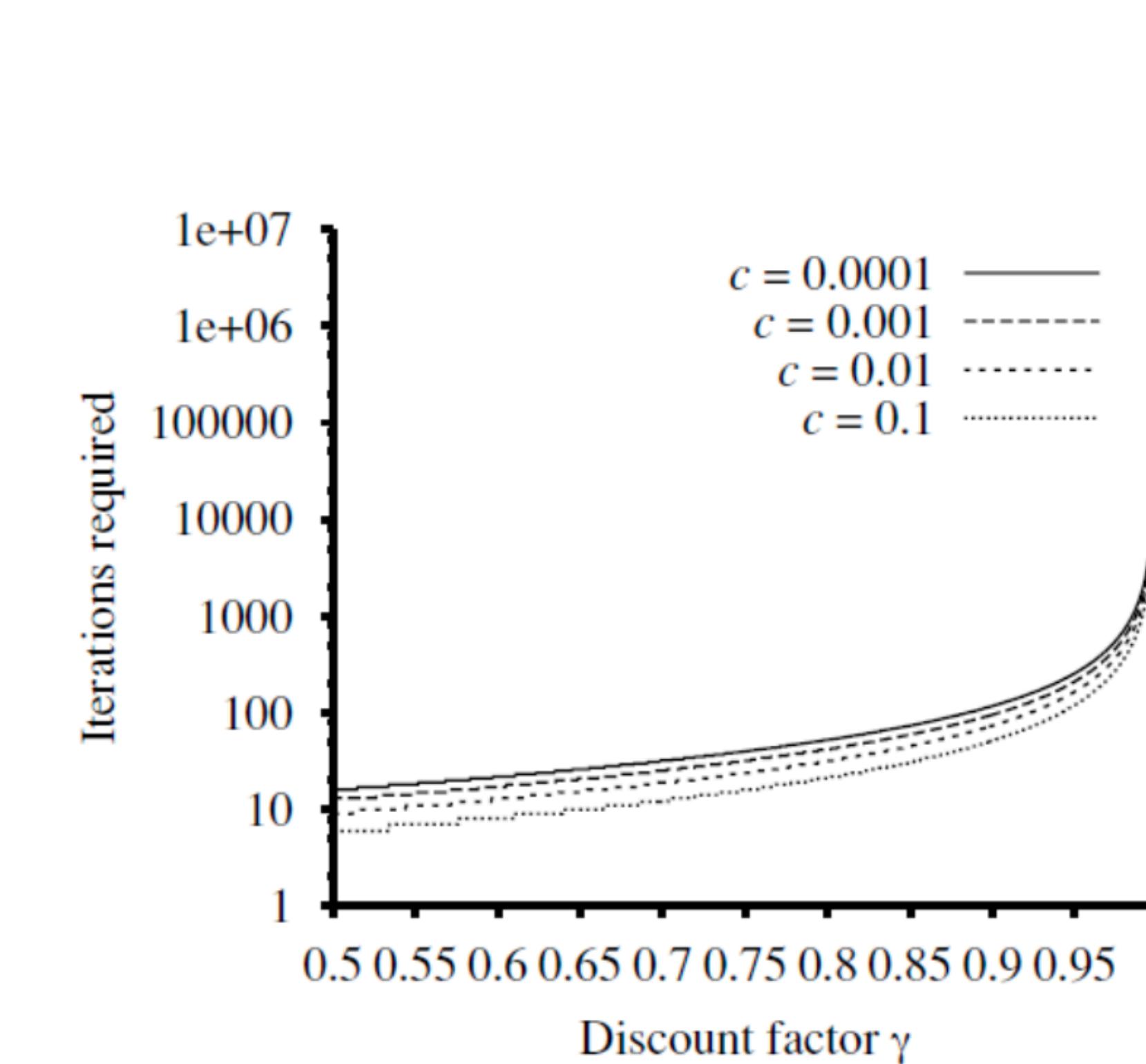
	1	2	3	4
3	.08	.56	.832	Butter
2	-.08		.464	Fire
1	-.08	-.08	-.08	-.08

# R

## Example



$$\gamma = 1$$



$$\epsilon \approx CR_{\max}$$

# R

# Solving MDPs

## Policy Iteration

- Start with some  $\pi_0$
- Repeat until convergence:

- Policy Evaluation  $V_{k+1}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_k(s))V_k(s')$
- Policy Improvement  $\pi_{k+1}(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a)V_{k+1}(s')$

$$\gamma = 1$$

$$O(n^3)$$

# R

# Solving MDPs

## Policy Iteration

- Start with some  $\pi_0$ : All up  $V([4,3]) = 1, V([4,2]) = -1, V(\cdot) = \underline{\underline{0}}$
- Compute values:

$$\begin{aligned} V_1([3,3]) &= -.04 + .1(0) + .8(0) + .1(1) = .06 \\ V_1([3,2]) &= -.04 + .1(0) + .8(0) + .1(-1) = -.14 \\ V_1([2,3]) &= -.04 + .1(0) + .8(0) + .1(0) = -.04 \end{aligned}$$

- Compute updated policy:  $\sum_{s'} P(s'|s, a)V_{k+1}(s')$

$$\begin{aligned} \text{For [3,3]:} \\ \text{Right} & \quad .1(.06) + .8(1) + .1(-.14) = \boxed{.792} \\ \text{Up} & \quad .1(-.14) + .8(.96) + .1(1) = .144 \\ \text{Left} & \quad .1(-.14) + .8(-.04) + .1(.06) = .014 \\ \text{Down} & \quad .1(-.04) + .8(-.14) + .1(1) = -.1 \end{aligned}$$

$$\Rightarrow \pi_1([3,3]) = \text{Right}$$

		+1	
		-1	