# Dynamic Programming (DP)

Mohammad **Dehghani**

Mechanical and Industrial Engineering Department

Northeastern University

Feb-20

- **Markov Decision Process**
  - Captures these two aspects of real-world problems:
    - » Involves an associative aspect—choosing different actions in different situations
    - » Actions influence not just immediate rewards, but also subsequent situations

- **MDP Solvers based on:**
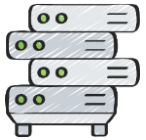  - the MDP model
  - the use if approximation methods

| | | Exact Representation | |
|---|---|---|---|
| | | Yes | No |
| **Known Model** | Yes | Dynamic Programming ❶ | Approximate DP ❷ |
| | No | Reinforcement Learning ❸ | RL with Function Approximation ❹ |

## What is Dynamic Programming?

- **_Dynamic_** sequential or temporal component to the problem
- **_Programming_** optimizing a "program", i.e. a policy
  - » c.f. linear programming

## DP is a method for solving complex problems, by breaking them down into a series of overlapping sub-problems

- Solve the subproblems
- Combine solutions to subproblems

## DP uses the _memorization_ technique.

- The technique of storing solutions to subproblems instead of recomputing them is called memorization.
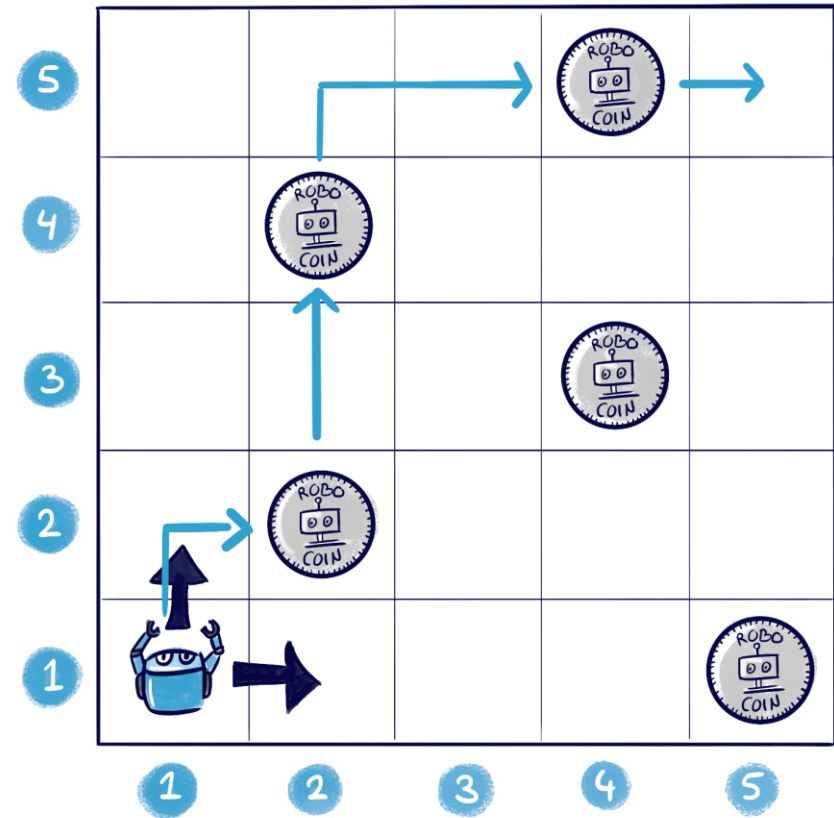
# Richard Bellman

**1920-1984**

**American Applied Mathematician**

- Introduced *Dynamic Programming (DP)* as a method for solving a complex problem by breaking it down into a collection of simpler subproblems, solving each of those subproblems just once, and storing their solutions.

- Bellman also introduced the *curse of dimensionality* which is an expression coined by Bellman to describe the problem caused by the exponential increase in volume associated with adding extra dimensions to a space.

# Dynamic Programming



DYNAMIC PROGRAMMING

**Dynamic Programming is a suitable method for problems with:**

**①** **Optimal substructure**

- Principle of optimality applies
- Optimal solution can be decomposed into subproblems
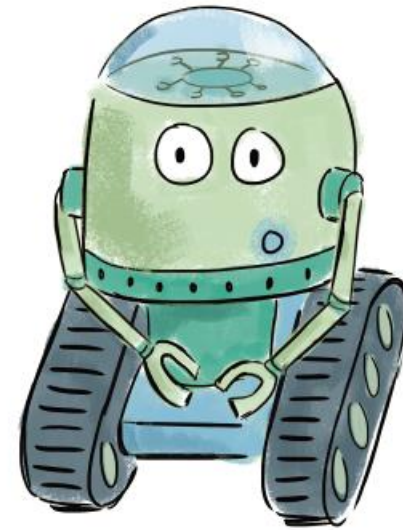
**②** **Overlapping subproblems**

- Subproblems recur many times
- Solutions can be cached and reused

- **Markov decision processes satisfy both properties**
  - Bellman equation gives recursive decomposition
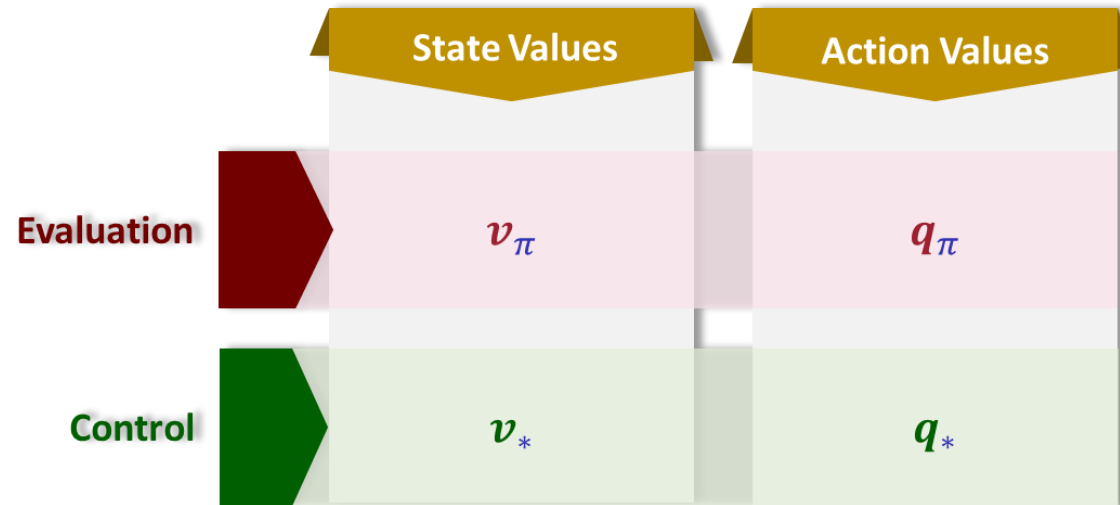  - Value function stores and reuses solutions

# Outline

▶ **The goal of this lecture is to develop new ways of calculating an optimal policy**

① **Policy evaluation**
- calculating value function for an arbitrary policy

② **Policy Improvement**
- obtain the new policy based on the new selected action

③ **Policy Iteration**
- calculating an optimal value function (and policy) by iteratively calculating value function and then improving policy

▶ **Discuss efficiency and utility of DP**

# *Policy Evaluation*
## vs.
## Control

# Policy Evaluation vs. Control

- **Policy evaluation**
  - The task of determining the value function for a for an arbitrary policy $\pi$

- **Control:**
  - The task of finding a policy to obtain as much reward as possible.



| | State Values | Action Values |
|---|---|---|
| **Evaluation** | $v_\pi$ | $q_\pi$ |
| **Control** | $v_*$ | $q_*$ |

▶ **The ultimate goal of is finding a policy which maximizes the value function (control)**

▶ *Policy Evaluation* **is a medium to reach** *Control*

$$\pi \rightarrow v_\pi$$

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

**According to policy $\pi$:**

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

- **Bellman equation reduces the problem of finding state-values $v_\pi(s)$ to a system of linear equations**

**Bellman equation**

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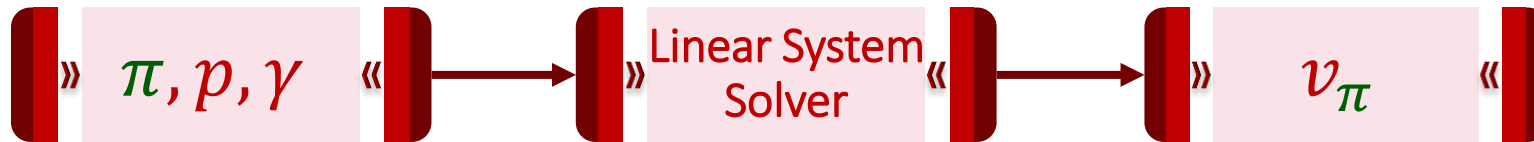
$$\pi \rightarrow v_\pi$$

🏴 **Simple MDPs:** Linear Programing

| » $\pi, p, \gamma$ « | → | » Linear System Solver « | → | » $v_\pi$ « |

✅ **General MDP:** Dynamic Programing

| » $\pi, p, \gamma$ « | → | » Dynamic Programming « | → | » $v_\pi$ « |

- **Control** corresponds to the operation of RL's policy-improvement algorithms.



$$\pi_2(s) > \pi_1(s)$$

$\pi_2$ strictly better than $\pi_1$

- **Dynamic programming assumes full knowledge of the MDP**

- **It is used for planning in an MDP**

  - For ***Prediction***:
    - » Input: MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma \rangle$ and policy $\pi$
    - » Output: value function $\boldsymbol{v_\pi}$

  - For ***Control***:
    - » Input: MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma \rangle$
    - » Output:
      - Optimal value function $\boldsymbol{v_*}$
      - Optimal policy $\boldsymbol{\pi_*}$

- The key idea of DP, and of RL generally, is the use of **value functions** to organize and structure the search for good policies.
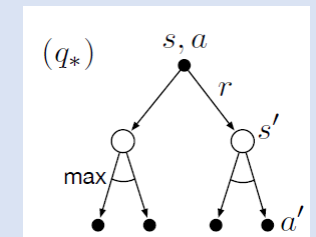
## Definition

- The **value of a state** under an optimal policy must equal the expected return for the **best** action from that state

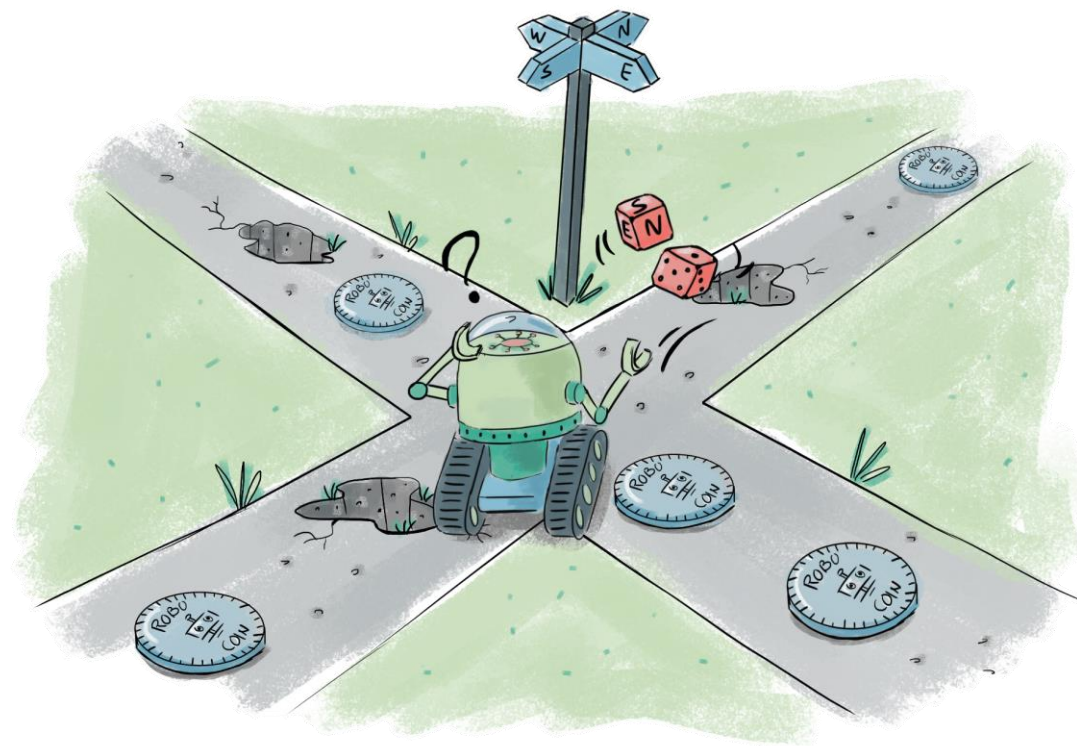$$v_*(s) = \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma \, v_*(s')], \text{ for all } s \in \mathcal{S}$$



- The **state-action pair value** under an optimal policy must equal the expected return for the **best** actions from the next state

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

# Policy Evaluation



POLICY EVALUATION

# Policy Evaluation

## Update Rule for Bellman Equation (state values)

**Consider a sequence of approximate value functions $v_0, v_1, v_2, \cdots$, each mapping $\mathcal{S}^+$ to $\mathbb{R}$ (the real numbers).**

- Using the Bellman equation, each successive approximation for state value is obtained by using the update rule as follows:

$$v_{k+1}(s) \doteq \mathbb{E}_\pi[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')], \text{ for } all \ s \in \mathcal{S}$$

where $v_k = v_\pi$ is a fixed point for this update rule because the Bellman equation for $v_\pi$ assures us of equality in this case

- The initial approximation, $v_0$, is chosen arbitrarily.
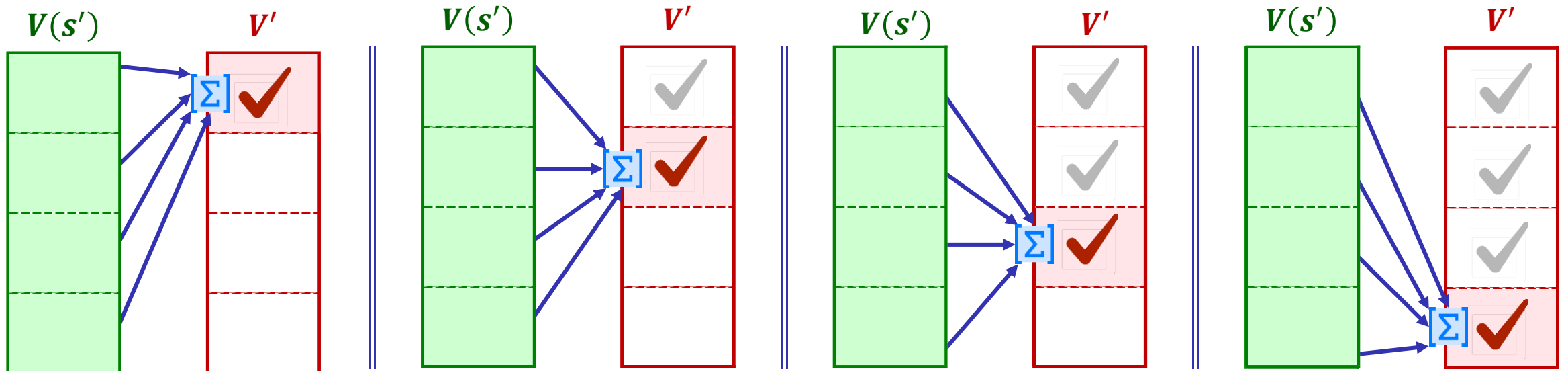- The terminal state, if any, must be given value 0

## Iterative policy evaluation

### Update Rule for Bellman Equation (state values)

$$v_{k+1}(s) \doteq \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')], \text{ for } all \ s \in \mathcal{S}$$

- $v_{k+1}(s)$ is the estimate of value function, which iteratively gets updated sequentially

- The sequence $\{v_k\}$ can be shown in general to converge to $v_\pi$ as $k \to \infty$

- The existence and uniqueness of $v_\pi$ are guaranteed as long as
  - $K \to \infty$
  - eventual termination is guaranteed from all states under the policy $\pi$

$$\lim_{k \to \infty} v_k = v_\pi, \qquad \forall \ v_0$$

## Update Rule for Bellman Equation (state values)

$$V'(s) \doteq \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')], \text{ for all } s \in \mathcal{S}$$
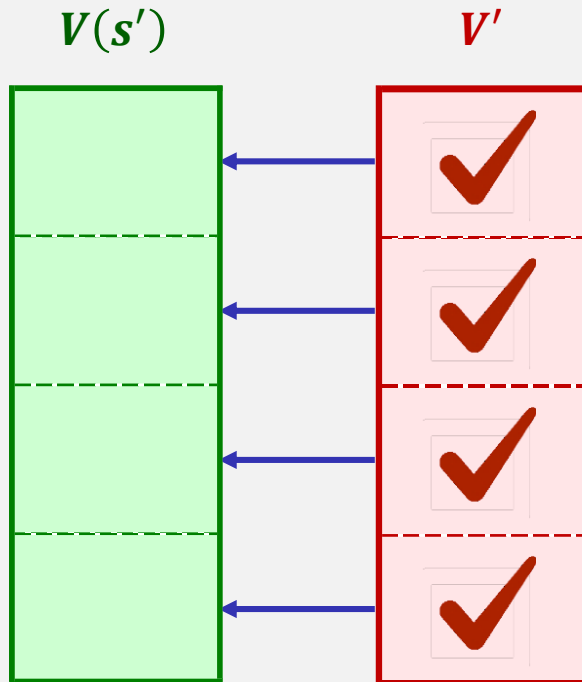
**Current state updated value**

**successive states (old values)**

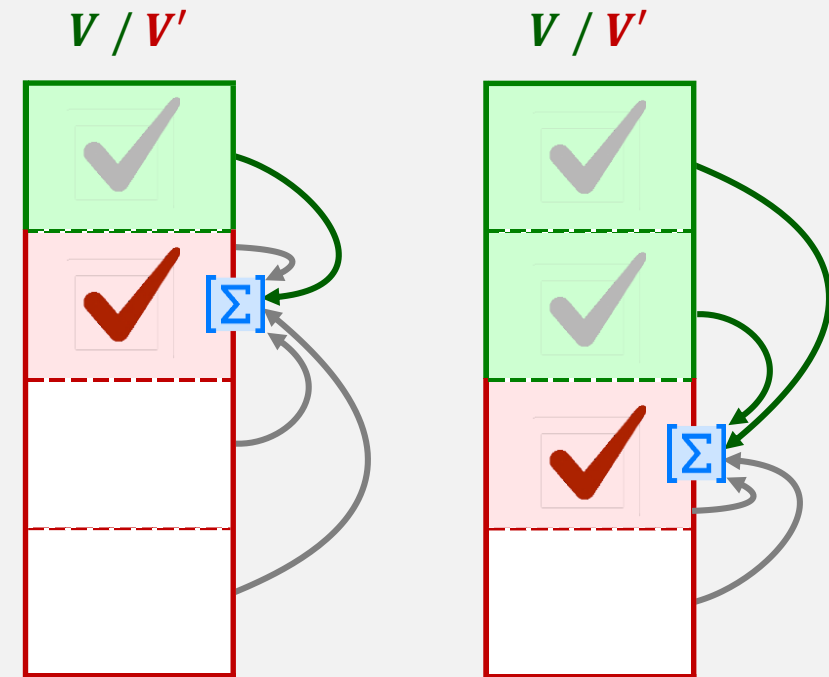## Two Array

$V(s')$       $V'$

## One Array
### In-Place Algorithm

$V$ / $V'$       $V$ / $V'$



- **In-Place Algorithm usually converges faster than the two-array version**
  - it uses new data as soon as they are available

Iterative policy evaluation.

**Iterative Policy Evaluation, for estimating $V \approx v_\pi$**

Input $\pi$, the policy to be evaluated
Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

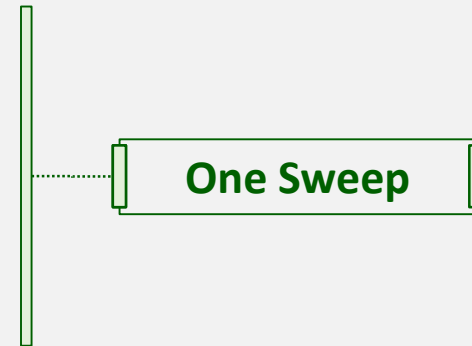Loop:
    $\Delta \leftarrow 0$
    Loop for each $s \in \mathcal{S}$:
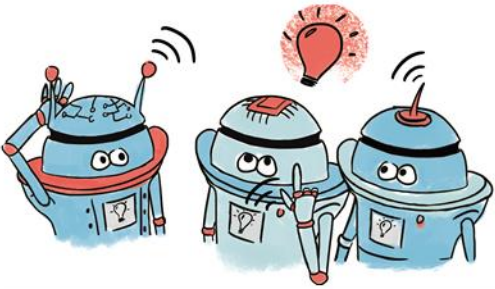        $v \leftarrow V(s)$
        $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

**One Sweep**

🔵 **Recall the Gridworld Example**

🔵 Use the evaluate iteration to calculate the state value of x for its 1st and 2nd iteration

   » $v_1(X) = ?$

   » $v_2(X) = ?$

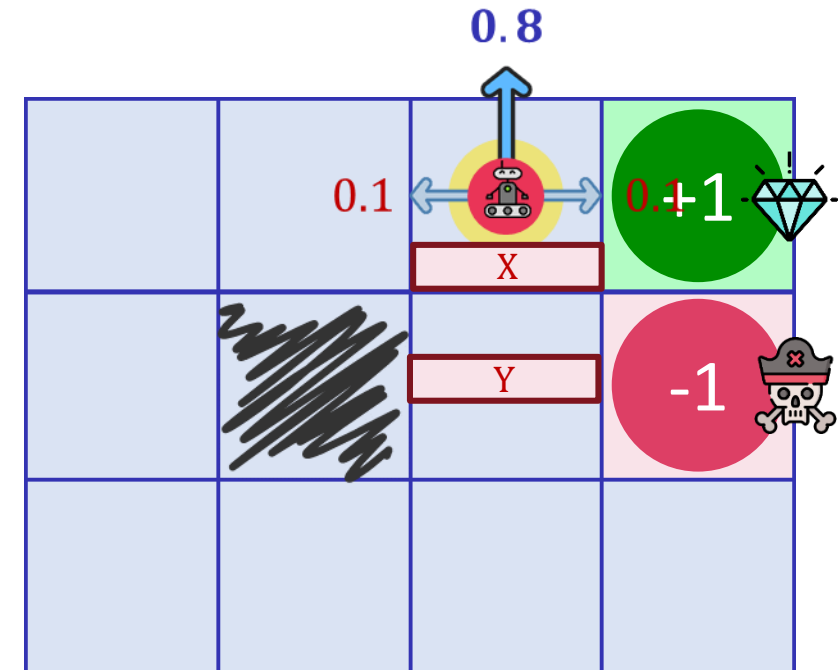🔴 **Parameters:**

   » $\gamma = 0.5$

   » $R(S) = -0.04$

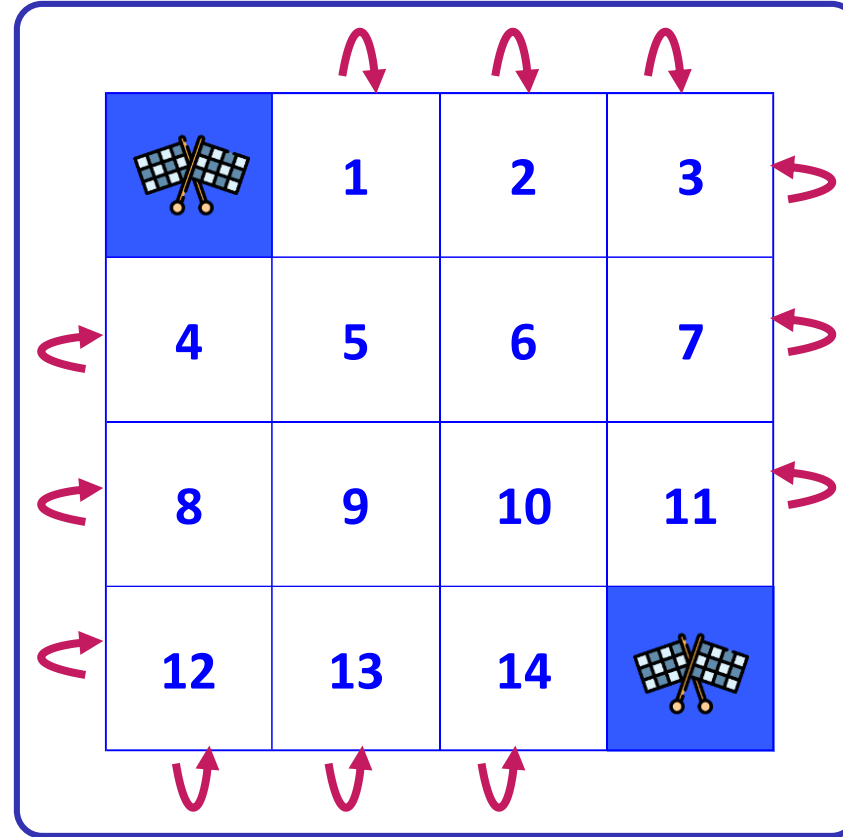   » $v_0(s) = 0$

🔴 **Policy**

   » Random Policy

$$\pi(n|\cdot) = \pi(e|\cdot) = \pi(w|\cdot) = \pi(s|\cdot) = 0.25$$

$$R_t = -1$$

$$\gamma = 1$$

**Reward**

**States**

**Actions**

$V$

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$$» \quad k = 1 \quad «$$

$$v_1(1) = \uparrow 0.25 \times (-1 + 0) + \\ \downarrow 0.25 \times (-1 + 0) + \\ \leftarrow 0.25 \times (-1 + 0) + \\ \rightarrow 0.25 \times (-1 + 0)$$

$$v_1(1) = -1$$

$V'$

| | | | |
|---|---|---|---|
| | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | |

$$v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma\, v_k(s')]$$

updated state value

$1$ $-1$

old values of the successor states

$V$

| 0 | -1 | -1 | -1 |
|---|----|----|----|
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 0 |

$$» \quad k = 2 \quad «$$

$v_2(1) = 0.25 \times (-1 + 1 \times -1) +$
$\qquad 0.25 \times (-1 + 1 \times -1) +$
$\qquad 0.25 \times (-1 + 1 \times 0) +$
$\qquad 0.25 \times (-1 + 1 \times -1) +$

$v_2(1) = -1.75$

$V'$

|  | -1.75 | -2 | -2 |
|---|----|----|----|
| -1.75 | -2 | -2 | -2 |
| -2 | -2 | -2 | -1.75 |
| -2 | -2 | -1.75 |  |

$$v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma \, v_k(s')]$$

1     −1

updated
state value

old values of the
successor states

## Example: $4 \times 4$ Gridworld

| » $k = 0$ « | » $k = 1$ « | » $k = 3$ « | » $k = 10$ « | » $k = \infty$ « |
|---|---|---|---|---|

$v_k(S)$

**$k = 0$**

|  | 0.00 | 0.00 | 0.00 |
|---|---|---|---|
| 0.00 | 0.00 | 0.00 | 0.00 |
| 0.00 | 0.00 | 0.00 | 0.00 |
| 0.00 | 0.00 | 0.00 |  |

**$k = 1$**

|  | -1.75 | -2.00 | -2.00 |
|---|---|---|---|
| -1.75 | -2.00 | -2.00 | -2.00 |
| -2.00 | -2.00 | -2.00 | -1.75 |
| -2.00 | -2.00 | -1.75 |  |

**$k = 3$**

|  | -2.44 | -2.94 | -3.00 |
|---|---|---|---|
| -2.44 | -2.88 | -3.00 | -2.94 |
| -2.94 | -3.00 | -2.88 | -2.44 |
| -3.00 | -2.94 | -2.44 |  |

**$k = 10$**

|  | -6.14 | -8.35 | -8.97 |
|---|---|---|---|
| -6.14 | -7.74 | -8.43 | -8.35 |
| -8.35 | -8.43 | -7.74 | -6.14 |
| -8.97 | -8.35 | -6.14 |  |

**$k = \infty$**

|  | -14.00 | -20.00 | -22.00 |
|---|---|---|---|
| -14.00 | -18.00 | -20.00 | -20.00 |
| -20.00 | -20.00 | -18.00 | -14.00 |
| -22.00 | -20.00 | -14.00 |  |

**Greedy Policy** w.r.t $v_k$



**Random Policy**
equiprobable

**Optimal Policy**

**Recall the Gridworld Example**

$$» \quad k = \infty \quad «$$

| | -14.0 | -20.0 | -22.0 |
|---|---|---|---|
| -14.0 | -18.0 | -20.0 | -20.0 |
| -20.0 | -20.0 | -18.0 | -14.0 |
| -22.0 | -20.0 | -14.0 | |

**1**  Why the model converges even with $\gamma = 1.0$

**2**  What is the simple interpretation of the values of states when the policy evaluation converges?

# Policy Improvement



POLICY IMPROVEMENT

# Policy Improvement

***Policy Evaluation*** helps us to compute value function for a policy $v_\pi(s)$

- $v_\pi(s)$: indicates know how good it is to follow the current policy from $s$

- $\pi \rightarrow \pi'$: results for policy evaluation can help develop better policy

***Policy Change*** *is applied after policy evaluation is performed*

1. $a \neq \pi(s)$: change in the policy at a single state to a particular action

2. $\pi'$: obtain the new policy based on the new selected action

3. $v_{\pi'}(s)$: evaluate the new changed policy

## Policy improvement theorem

- **Let $\pi$ and $\pi'$ be any pair of deterministic policies such that**

$$q_\pi\big(s, \pi'(s)\big) \geq q_\pi\big(s, \pi(s)\big), \text{for all } s \in \mathcal{S},$$

New Value          Old value

▶ then the policy $\pi'$ must be as good as, or better than, $\pi$.

$$\pi' \geq \pi$$

▶ The new policy must obtain greater or equal expected return from all states $s \in \mathcal{S}$ :

$$v_{\pi'}(s) \geq v_\pi(s), \text{for all } s \in \mathcal{S}$$

- **The new policy is a *strict* improvement over $\pi$, if and only if:**

$$q_\pi\big(s, \pi'(s)\big) > q_\pi\big(s, \pi(s)\big) \text{ for at least one } s \in \mathcal{S} \quad \rightarrow \quad \pi' > \pi$$

## *How to select actions to improve the policy?*

> **Greedy policy**
>
> **Greedy policy** considers changes at all states and to all possible actions, by selects at each state the action that appears best according to $q_\pi(s, a)$
>
> $$\pi'(s) \doteq \underset{a}{\arg\max}\, q_\pi(s, a)$$
>
> $$\doteq \underset{a}{\arg\max} \sum_{s',r} p(s', r | s, a)[r + \gamma v_\pi(s')],$$

- By construction, the greedy policy meets the conditions of the policy improvement theorem
- We know that the new policy is as good as, or better than, the original policy.

## Policy Improvement

**The process of making a new policy that improves on an original policy, by making it greedy with respect to the value function of the original policy, is called policy improvement.**

$$v_{\pi'}(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_{\pi'}(S_{t+1})|S_t = s, A_t = a]$$
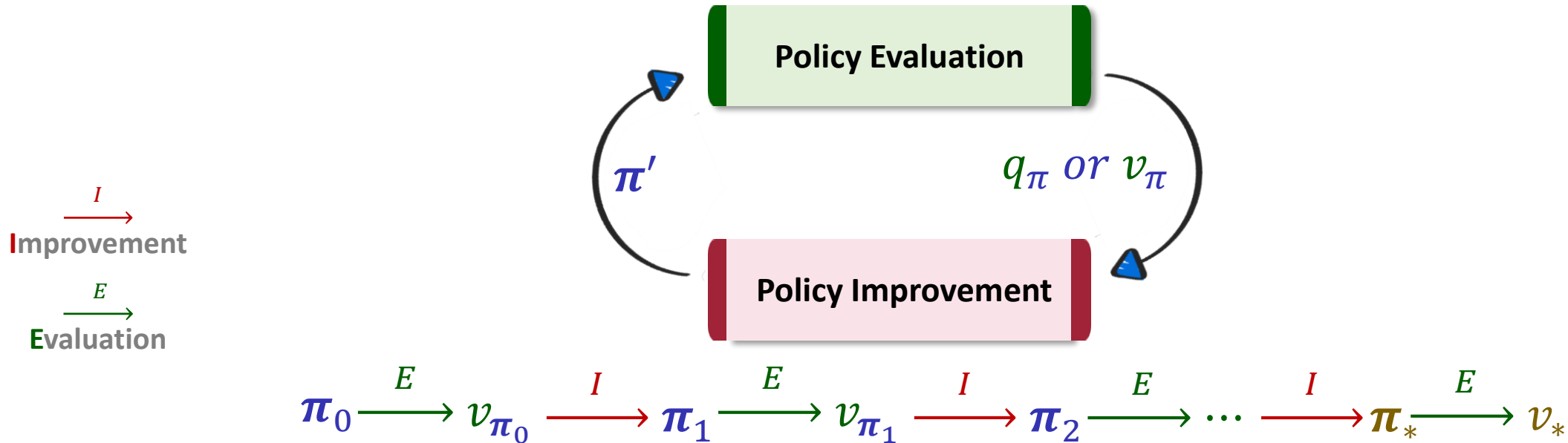
$$v_{\pi'}(s) = \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_{\pi'}(s')],$$

- Policy improvement is same as the Bellman optimality equation, and therefore:
  - ☑ $v_{\pi'}(s)$ must be $v_*$
  - ☑ **Both $\pi'$ and $\pi$** must be optimal policies
- Policy improvement thus must give us a strictly better policy except when the original policy is already optimal.

# Policy Iteration

# Policy Iteration

- **Policy iteration:** is the process combining policy evaluation, updating the value of the policy, and policy improvement, obtaining the best policy available

$$\xrightarrow{\ I\ }$$
**I**mprovement

$$\xrightarrow{\ E\ }$$
**E**valuation



$$\pi_0 \xrightarrow{\ E\ } v_{\pi_0} \xrightarrow{\ I\ } \pi_1 \xrightarrow{\ E\ } v_{\pi_1} \xrightarrow{\ I\ } \pi_2 \xrightarrow{\ E\ } \cdots \xrightarrow{\ I\ } \pi_* \xrightarrow{\ E\ } v_*$$

- Each policy is guaranteed to be a strict improvement over the previous one
- Because a finite MDP:
  - » has only a finite number of policies
  - » this process must converge to an optimal policy and optimal value function in a finite number of iterations

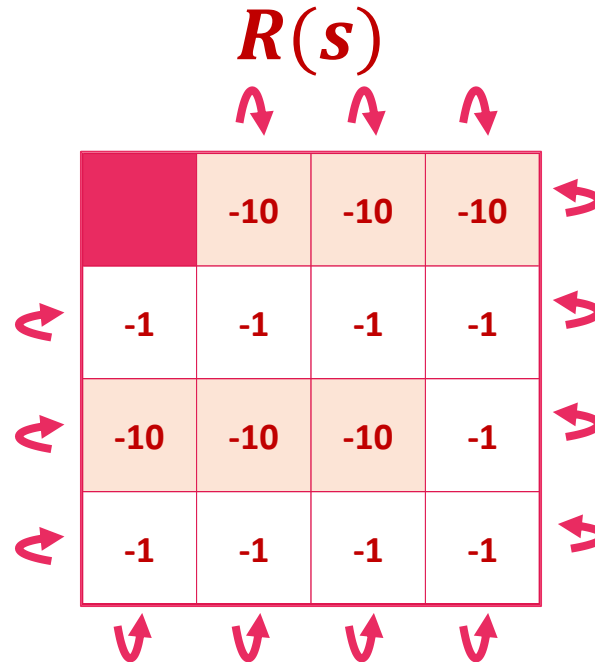# Policy Iteration

- In Small Gridworld improved policy was optimal, $\pi' =$

- In general, need more iterations of improvement / evaluation

- But this process of policy iteration always converges to $\pi_*$

$$R(s)$$

| | -10 | -10 | -10 |
|---|---|---|---|
| -1 | -1 | -1 | -1 |
| -10 | -10 | -10 | -1 |
| -1 | -1 | -1 | -1 |

**Excel Lab**

The Gridworld Example: New Version

# Policy Iteration

**Old Policy**

| | | | |
|---|---|---|---|
| | ← | ← | ← |
| ↑ | ← | ← | ← |
| ↑ | ↑ | ↑ | ↑ |
| ↑ | ← | ← | ← |

**State Vlaues**

| | -1 | -11 | -21 |
|---|---|---|---|
| -1 | -2 | -3 | -4 |
| -2 | -3 | -4 | -5 |
| -12 | -13 | -14 | -15 |

**New Policy**

| | ← | ↓ | ↓ |
|---|---|---|---|
| ↑ | ← | ← | ← |
| ↑ | ↑ | ↑ | ↑ |
| ↑ | ←↑ | ←↑ | ↑ |

# Policy Iteration

- **We repeat iterations until the policy becomes stable**



| Old Policy | Optimal Policy | Optimal State Vlaues |
|:---:|:---:|:---:|

**Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Loop:
   $\quad\quad \Delta \leftarrow 0$
   $\quad\quad$ Loop for each $s \in \mathcal{S}$:
   $\quad\quad\quad v \leftarrow V(s)$
   $\quad\quad\quad V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s))\big[r + \gamma V(s')\big]$
   $\quad\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
   $\quad\quad old\text{-}action \leftarrow \pi(s)$
   $\quad\quad \pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
   $\quad\quad$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

PAIR, THINK, SHARE

● **Policy Iteration pseudocode**

The policy iteration algorithm on page 80 has a subtle bug in that it may never terminate if the policy continually switches between two or more policies that are equally good.

● **Describe qualitatively how we can modify the pseudocode so that convergence is guaranteed?**

**Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Loop:
   $\quad \Delta \leftarrow 0$
   $\quad$ Loop for each $s \in \mathcal{S}$:
   $\quad\quad v \leftarrow V(s)$
   $\quad\quad V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s))[r + \gamma V(s')]$
   $\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
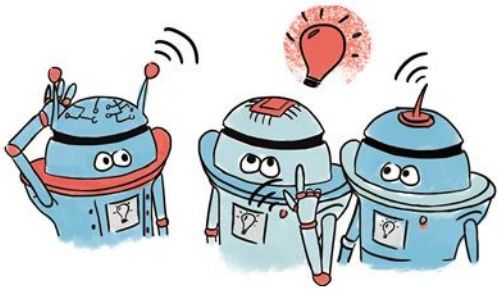   $policy\text{-}stable \leftarrow true$
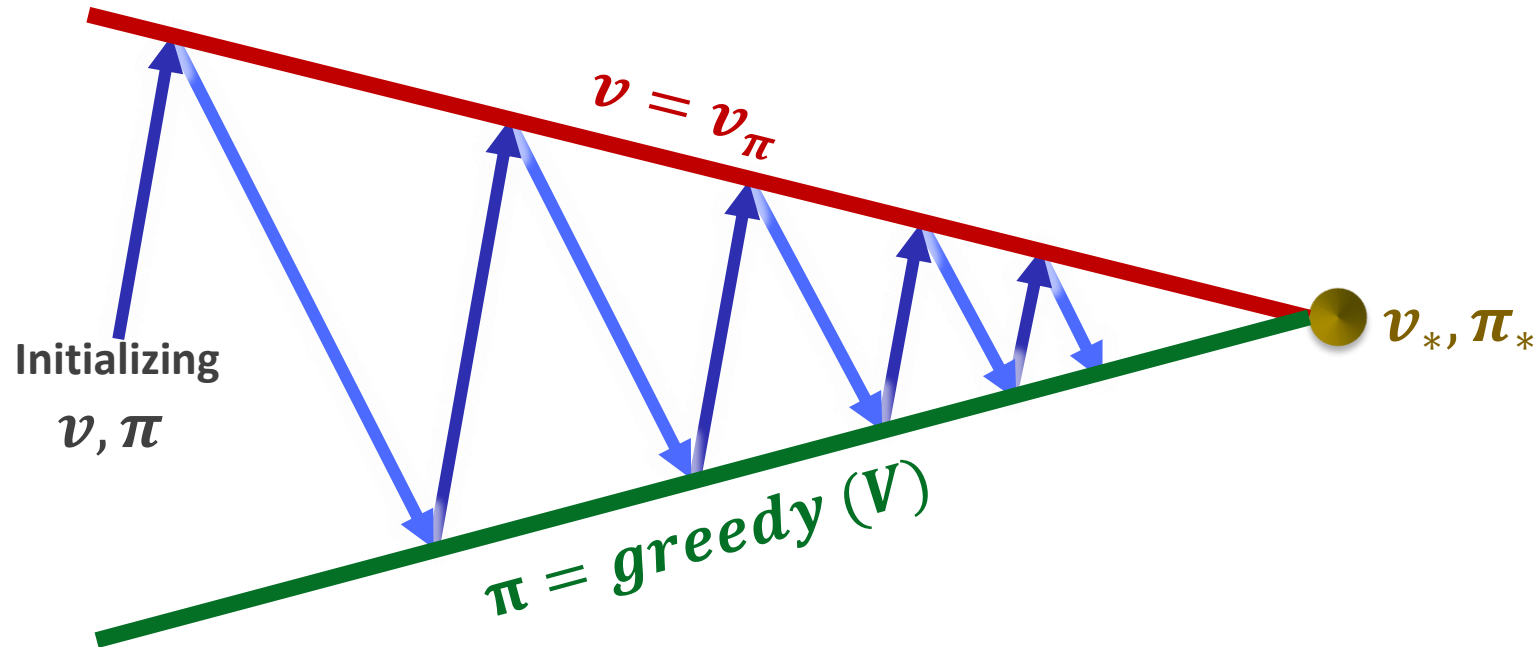   For each $s \in \mathcal{S}$:
   $\quad old\text{-}action \leftarrow \pi(s)$
   $\quad \pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
   $\quad$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

$$v = v_\pi$$

$$v_*, \pi_*$$

**Initializing**

$$v, \pi$$

$$\pi = greedy\,(V)$$

evaluation

$$V \rightsquigarrow v_\pi$$

$$\pi$$      $$V$$

$$\pi \rightsquigarrow greedy(V)$$
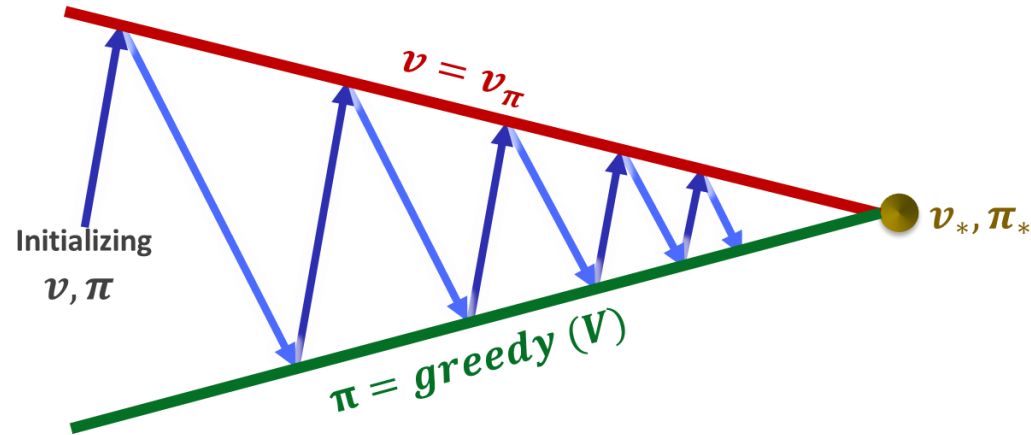
improvement

$$\pi_* \rightleftharpoons v_*$$

- Any policy evaluates to a unique value function which can be greedified to produce a better policy
- That in turn evaluates to a value function which can in turn be greedified
- Each policy is *strictly better* than the previous, until *eventually both are optimal*
- The dance converges in a finite number of steps, usually very few

- 🔴 **Competing:**
  They compete in the sense that they pull in opposing directions. Making the policy greedy with respect to the value function typically makes the value function incorrect for the changed policy, and making the value function consistent with the policy typically causes that policy no longer to be greedy.
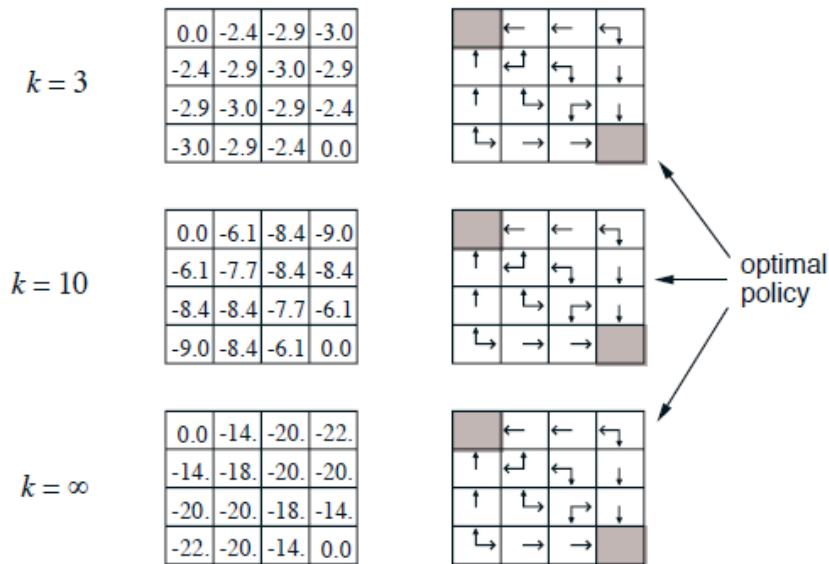
- 🔴 **Cooperating:**
  In the long run, however, these two processes interact to find a single joint solution

# Value Iteration

- **One drawback to policy iteration is that each of its iterations involves policy evaluation**
  - Iterative computation requiring multiple sweeps through the state set

- **Must we wait for exact convergence, or can we stop short of that?**



In this example, policy evaluation iterations beyond the first three have no effect on the corresponding greedy policy.

# Value Iteration

- **One drawback to policy iteration is that each of its iterations involves policy evaluation**
  - Iterative computation requiring multiple sweeps through the state set

- **In fact, the policy evaluation step of policy iteration can be truncated**
  - i.e. stop policy evaluation just one sweep (one update of each state).
  - This algorithm is called value iteration

**Update Rule for Bellman Equation (action values)**

$$v_{k+1}(s) \doteq \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma\, v_k(s')], \text{ for } all\ s \in \mathcal{S}$$

# Value Iteration

**Value Iteration, for estimating $\pi \approx \pi_*$**

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
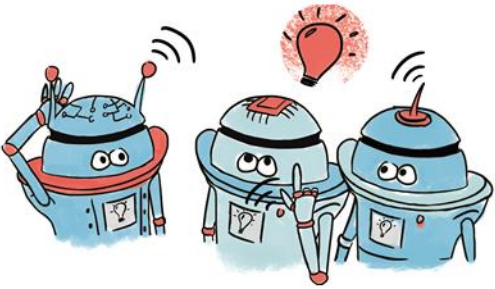Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
$\quad | \quad \Delta \leftarrow 0$
$\quad | \quad$ Loop for each $s \in \mathcal{S}$:
$\quad | \qquad v \leftarrow V(s)$
$\quad | \qquad V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)\left[r + \gamma V(s')\right]$
$\quad | \qquad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that
$\quad \pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)\left[r + \gamma V(s')\right]$

**Recall the Gridworld Example**

Use the evaluate iteration to calculate the state value of x for its 1st and 2nd iteration

» $v_1(X) =$?

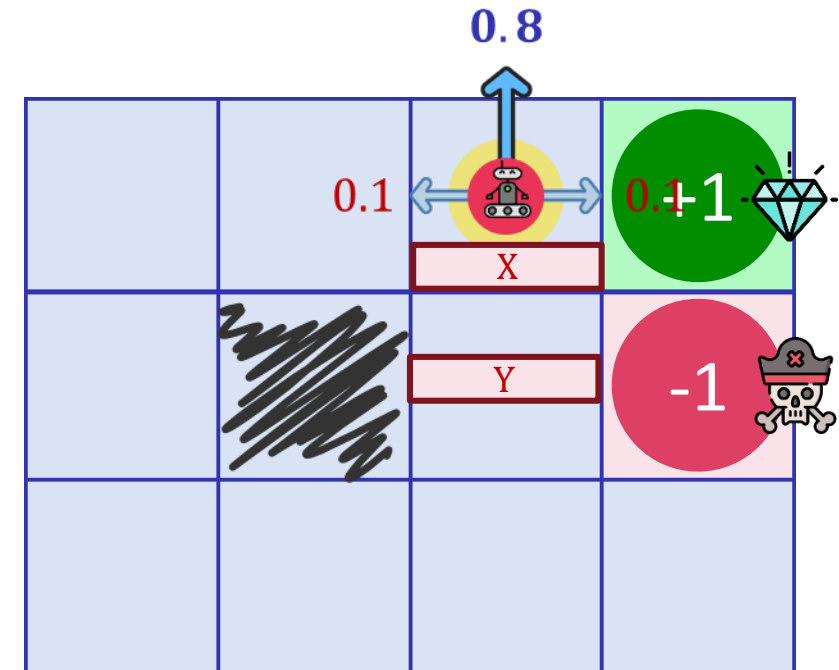» $v_2(X) =$?

**Parameters:**

» $\gamma = 0.5$

» $R(S) = -0.04$

» $v_0(s) = 0$

**Policy**
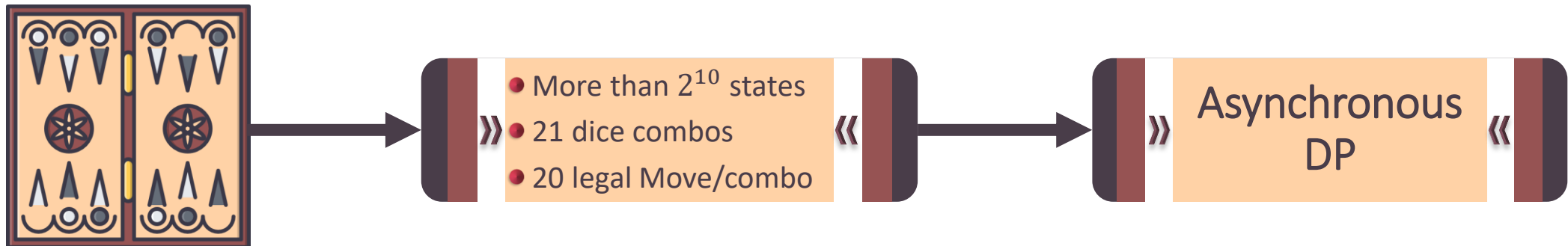
» Based on the Value Iteration?

# The curse of dimensionality

- **We can DP is appropriate when the optimal policy is polynomial in the number of states**



- **the number of states is often astronomical, e.g., often growing exponentially with the number of state** variables (what Bellman called "the curse of dimensionality")



- More than $2^{10}$ states
- 21 dice combos
- 20 legal Move/combo

Asynchronous DP

# Asynchronous DP

- DP methods described so far used synchronous backups
  - i.e. all states are backed up in parallel

- Asynchronous DP backs up states individually, in any order
  - For each selected state, apply the appropriate backup
  - Can significantly reduce computation
  - Guaranteed to converge if all states continue to be selected
    - To converge correctly, however, an asynchronous algorithm must continue to update the values of all the states: it can't ignore any state after some point in the computation.
  - take advantage of this flexibility by selecting the states to which we apply updates so as to improve the algorithm's rate of progress

# Summary

- **DP in RL follows these 3 steps:**
  - Policy Evaluation:
  - Policy Improvement
  - Policy Iteration

- **The improvement theorem is also valid for stochastic policy $\pi(a|s)$ cases**
  - In the general case, a stochastic policy $\boldsymbol{\pi}$ specifies probabilities, $\boldsymbol{\pi(a|s)}$, for taking each
  - action, a, in each state, s.
  - In particular, the policy improvement theorem carries through as stated for the stochastic case
  - If there are ties in policy improvement steps, each maximizing action can be given a portion of the probability of being selected in the new greedy policy

# References

- [1] R. S. Sutton and A. G. Barto, Reinforcement learning: An introduction. MIT press, 2018.