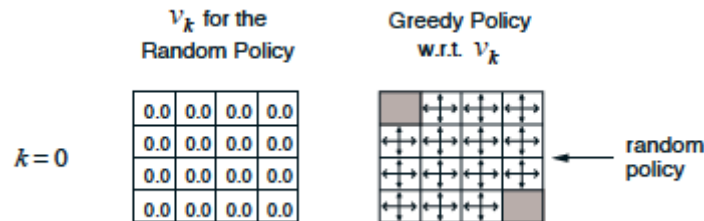




## Summary



First step of policy iteration in gridworld example (Sutton and Barto, 2017)

## Introduction

- In the **dynamic programming** setting, the agent has full knowledge of the MDP. (This is much easier than the **reinforcement learning** setting, where the agent initially knows nothing about how the environment decides state and reward and must learn entirely from interaction how to select actions.)

## An Iterative Method

- In order to obtain the state-value function  $v_\pi$  corresponding to a policy  $\pi$ , we need only solve the system of equations corresponding to the Bellman expectation equation for  $v_\pi$ .
- While it is possible to analytically solve the system, we will focus on an iterative solution approach.

## Iterative Policy Evaluation

- **Iterative policy evaluation** is an algorithm used in the dynamic programming setting to estimate the state-value function  $v_\pi$  corresponding to a policy  $\pi$ . In this approach, a Bellman update is applied to the value function estimate until the changes to the estimate are nearly imperceptible.



## Summary

```

Input: MDP, policy  $\pi$ , small positive number  $\theta$ 
Output:  $V \approx v_\pi$ 
Initialize  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^+$ )
repeat
   $\Delta \leftarrow 0$ 
  for  $s \in \mathcal{S}$  do
     $v \leftarrow V(s)$ 
     $V(s) \leftarrow \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r|s, a)(r + \gamma V(s'))$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  end
until  $\Delta < \theta$ ;
return  $V$ 

```

## Estimation of Action Values

- In the dynamic programming setting, it is possible to quickly obtain the action-value function  $q_\pi$  from the state-value function  $v_\pi$  with the equation:

$$q_\pi(s, a) = \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r|s, a)(r + \gamma v_\pi(s')).$$

### Estimation of Action Values

```

Input: state-value function  $V$ 
Output: action-value function  $Q$ 
for  $s \in \mathcal{S}$  do
  for  $a \in \mathcal{A}(s)$  do
     $Q(s, a) \leftarrow \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r|s, a)(r + \gamma V(s'))$ 
  end
end
return  $Q$ 

```

## Policy Improvement

- Policy improvement** takes an estimate  $V$  of the action-value function  $v_\pi$  corresponding to a policy  $\pi$ , and returns an improved (or equivalent) policy  $\pi'$ , where  $\pi' \geq \pi$ . The algorithm first constructs the action-value function estimate  $Q$ . Then, for each state  $s \in \mathcal{S}$ , you need only select the action  $a$  that maximizes  $Q(s, a)$ . In other words,  $\pi'(s) = \arg \max_{a \in \mathcal{A}(s)} Q(s, a)$  for all  $s \in \mathcal{S}$ .



## Summary

```

Input: MDP, value function  $V$ 
Output: policy  $\pi'$ 
for  $s \in \mathcal{S}$  do
    for  $a \in \mathcal{A}(s)$  do
         $Q(s, a) \leftarrow \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r | s, a)(r + \gamma V(s'))$ 
    end
     $\pi'(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} Q(s, a)$ 
end
return  $\pi'$ 

```

## Policy Iteration

- **Policy iteration** is an algorithm that can solve an MDP in the dynamic programming setting. It proceeds as a sequence of policy evaluation and improvement steps, and is guaranteed to converge to the optimal policy (for an arbitrary *finite* MDP).

### Policy Iteration

```

Input: MDP, small positive number  $\theta$ 
Output: policy  $\pi \approx \pi_*$ 
Initialize  $\pi$  arbitrarily (e.g.,  $\pi(a|s) = \frac{1}{|\mathcal{A}(s)|}$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ )
 $policy\_stable \leftarrow false$ 
repeat
     $V \leftarrow \text{Policy\_Evaluation}(\text{MDP}, \pi, \theta)$ 
     $\pi' \leftarrow \text{Policy\_Improvement}(\text{MDP}, V)$ 
    if  $\pi = \pi'$  then
         $policy\_stable \leftarrow true$ 
    end
     $\pi \leftarrow \pi'$ 
until  $policy\_stable = true$ ;
return  $\pi$ 

```

## Truncated Policy Iteration

- **Truncated policy iteration** is an algorithm used in the dynamic programming setting to estimate the state-value function  $v_\pi$  corresponding to a policy  $\pi$ . In this approach, the evaluation step is stopped after a fixed number of sweeps through the state space. We refer to the algorithm in the evaluation step as **truncated policy evaluation**.



## Summary

**Input:** MDP, policy  $\pi$ , value function  $V$ , positive integer  $max\_iterations$   
**Output:**  $V \approx v_\pi$  (if  $max\_iterations$  is large enough)  
 $counter \leftarrow 0$   
**while**  $counter < max\_iterations$  **do**  
  **for**  $s \in \mathcal{S}$  **do**  
     $V(s) \leftarrow \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{s', r \in \mathcal{R}} p(s', r|s, a)(r + \gamma V(s'))$   
  **end**  
   $counter \leftarrow counter + 1$   
**end**  
**return**  $V$

## Truncated Policy Iteration

**Input:** MDP, positive integer  $max\_iterations$ , small positive number  $\theta$   
**Output:** policy  $\pi \approx \pi_*$   
Initialize  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^+$ )  
Initialize  $\pi$  arbitrarily (e.g.,  $\pi(a|s) = \frac{1}{|\mathcal{A}(s)|}$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ )  
**repeat**  
   $\pi \leftarrow \text{Policy\_Improvement}(\text{MDP}, V)$   
   $V_{old} \leftarrow V$   
   $V \leftarrow \text{Truncated\_Policy\_Evaluation}(\text{MDP}, \pi, V, max\_iterations)$   
**until**  $\max_{s \in \mathcal{S}} |V(s) - V_{old}(s)| < \theta$ ;  
**return**  $\pi$

## Value Iteration

- **Value iteration** is an algorithm used in the dynamic programming setting to estimate the state-value function  $v_\pi$  corresponding to a policy  $\pi$ . In this approach, each sweep over the state space simultaneously performs policy evaluation and policy improvement.

## Value Iteration

**Input:** MDP, small positive number  $\theta$   
**Output:** policy  $\pi \approx \pi_*$   
Initialize  $V$  arbitrarily (e.g.,  $V(s) = 0$  for all  $s \in \mathcal{S}^+$ )  
**repeat**  
   $\Delta \leftarrow 0$   
  **for**  $s \in \mathcal{S}$  **do**  
     $v \leftarrow V(s)$   
     $V(s) \leftarrow \max_{a \in \mathcal{A}(s)} \sum_{s', r \in \mathcal{R}} p(s', r|s, a)(r + \gamma V(s'))$   
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$   
  **end**  
**until**  $\Delta < \theta$ ;  
 $\pi \leftarrow \text{Policy\_Improvement}(\text{MDP}, V)$   
**return**  $\pi$



Summary

---