

Policy Iteration Algorithm

Algorithm 1 Policy Iteration

```
1: Initialize a random policy  $\pi$ 
2: repeat
3:   Policy Evaluation:
4:   Initialize  $V(s) = 0$  for all states  $s$ 
5:   repeat
6:     for all states  $s$  do
7:        $V_{\text{new}}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ 
8:     end for
9:      $V \leftarrow V_{\text{new}}$ 
10:  until convergence
11:  Policy Improvement:
12:  Policy-stable  $\leftarrow$  true
13:  for all states  $s$  do
14:     $\text{old\_action} \leftarrow \pi(s)$ 
15:     $\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$ 
16:    if  $\text{old\_action} \neq \pi(s)$  then
17:      Policy-stable  $\leftarrow$  false
18:    end if
19:  end for
20: until Policy-stable
```
