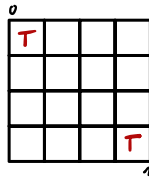


# 1. DP - Policy evaluation

Fixed  $\pi \rightarrow$  state-value function  $V^\pi$  for MDP:



A:  $\uparrow \downarrow \leftarrow \rightarrow$   
 stop on T:  $R = -1$   
 T:  $R = 0$

Stop:  $\Delta V < \phi$

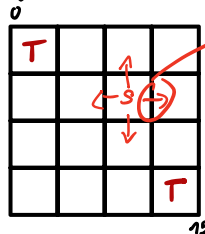
def policy-eval( $\pi, env, \gamma=1, \phi=0.0001$ )

$$V(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) [r(s,a,s') + \gamma V(s')]$$

for each  $s$ :  $N = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) [r + \gamma V(s')]$

## 2. Policy Iteration

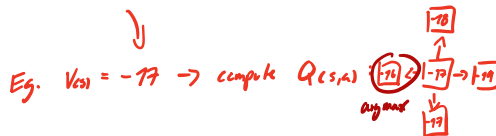
Same grid, actions...



Goal: Given  $V(s)$  improve the policy by acting greedily

For each state  $s$ : choose action that maximizes expected return

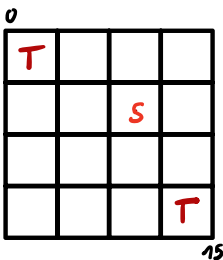
we use the  $V(s)$  from our env!



we update policy for this  $s$  to always go right

- Now run policy-eval again but with our improved policy  
 ~ and get  $V^*[s] = -2$

Same as 1 & 2:



## 3 Value Iteration

- updates the value of each state by considering the best possible action  $\rightarrow$  we extract optimal policy from it.

- after convergence the optimal policy is extracted  
 - requires full knowledge of env.

## 4. Gambler's problem

- Find best way to win \$100 before going broke
- $V[s]$  shows chance of reaching \$100 from capital  $s$
- using **value iteration** for best optimal policy
- optimal policy: often bet big (depends on  $p_h$ )

def one-step lookahead:

for each state in capital:

$$\text{Expected} = p_h \cdot (\text{reward\_win} + V[\text{win\_cap}]) + (1-p_h) \cdot (\text{reward\_lose} + V[\text{lose\_cap}])$$

def: val-iter-gambler:

init...  
 until delta < theta:  
 $V[s] = \max(\text{one-step lookahead})$

for states 1-99:

action-val = one-step lookahead

best action = max  $\rightarrow +1$

policy[s] =  $\rightarrow$