

RL POLICIES

BIOE 498/598 P1

Review

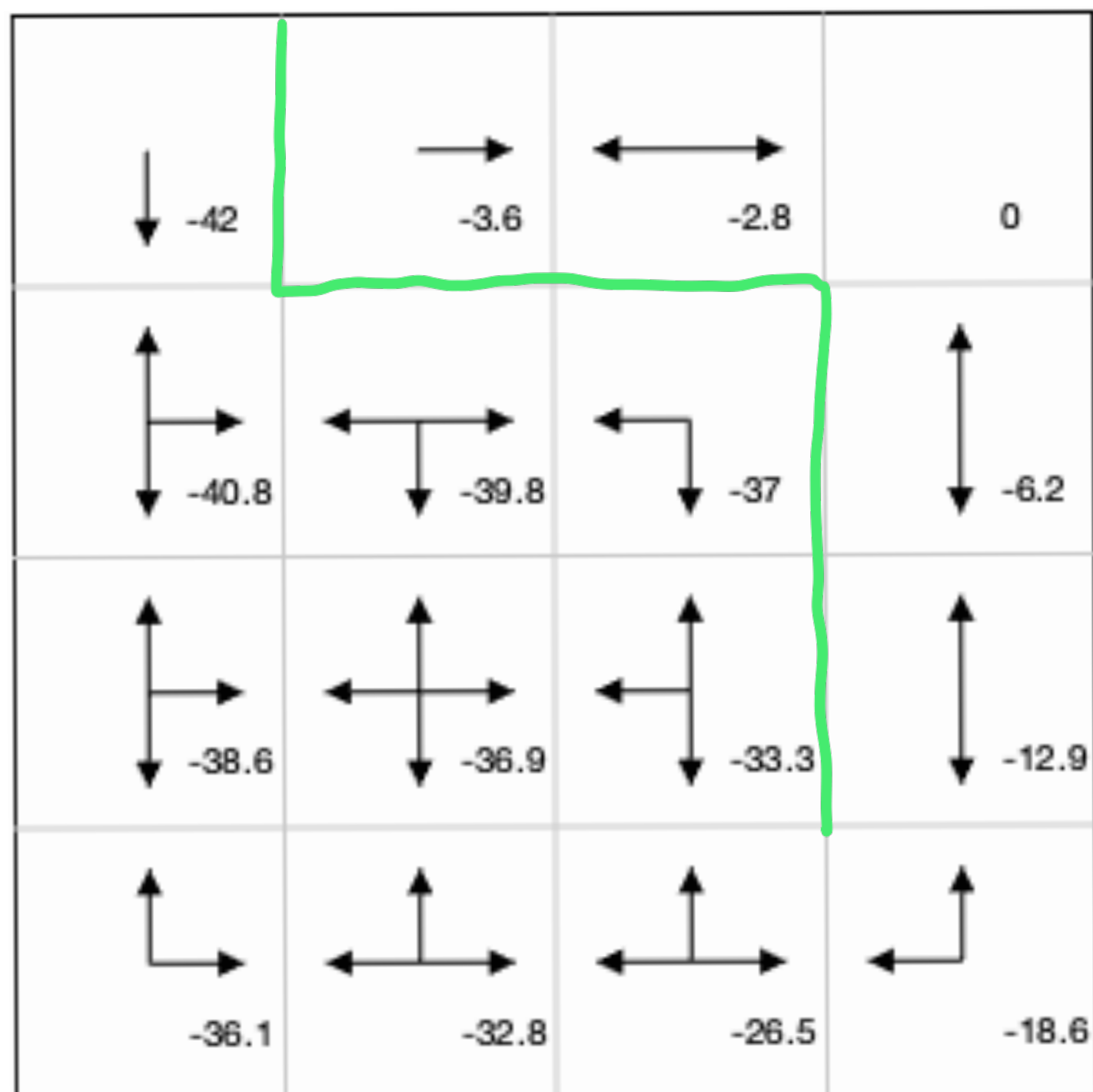
RL uses MDP framework



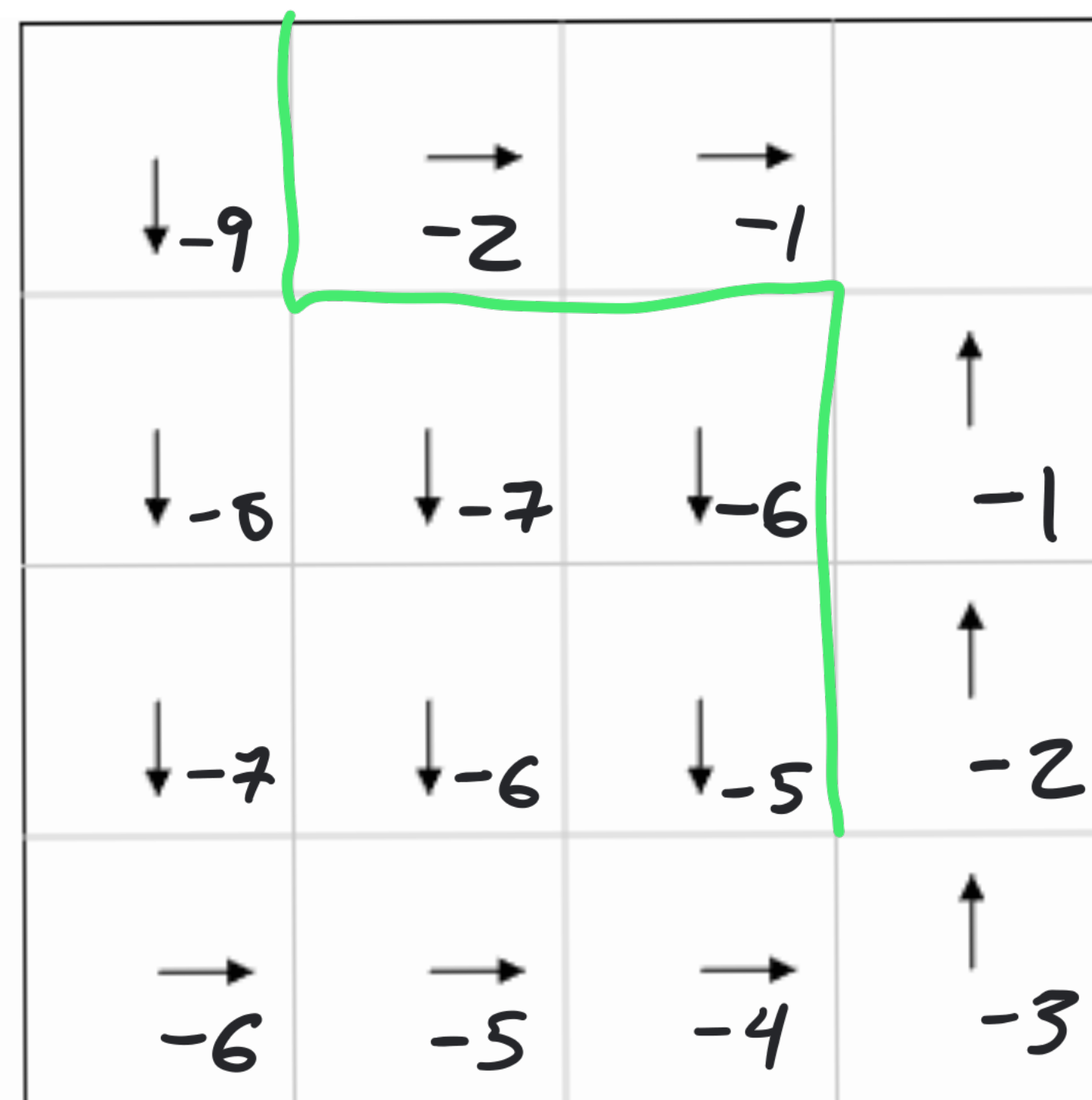
$$\max_a \left\{ r_1 + \sum_{i=2}^T r_i \right\} = \max_a \left\{ r_1 + V(s') \right\}$$

Policy π is a function that choose actions for every state.

$$a = \pi(s)$$



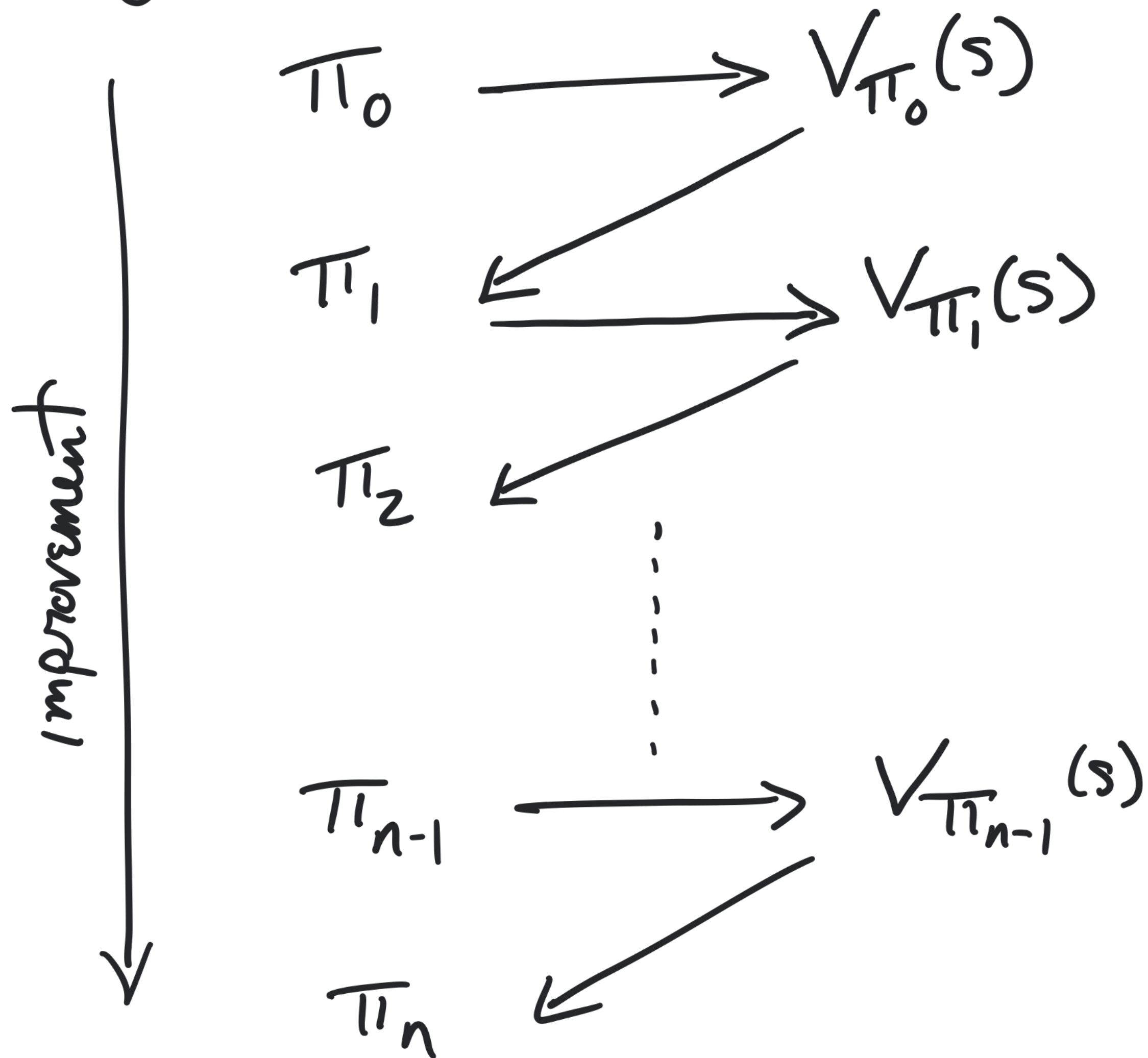
$V(s)$ under random moves



$V(s)$ new policy

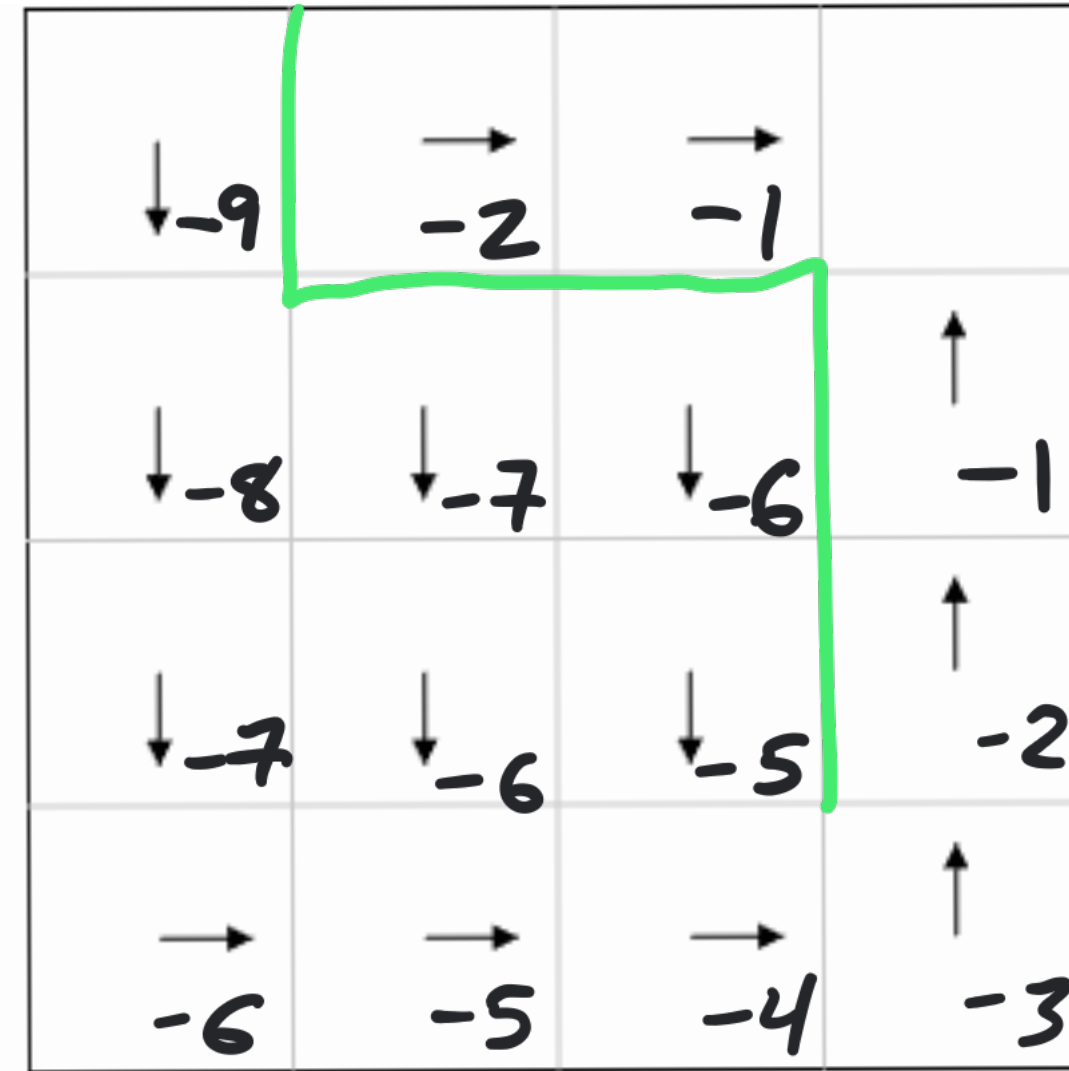
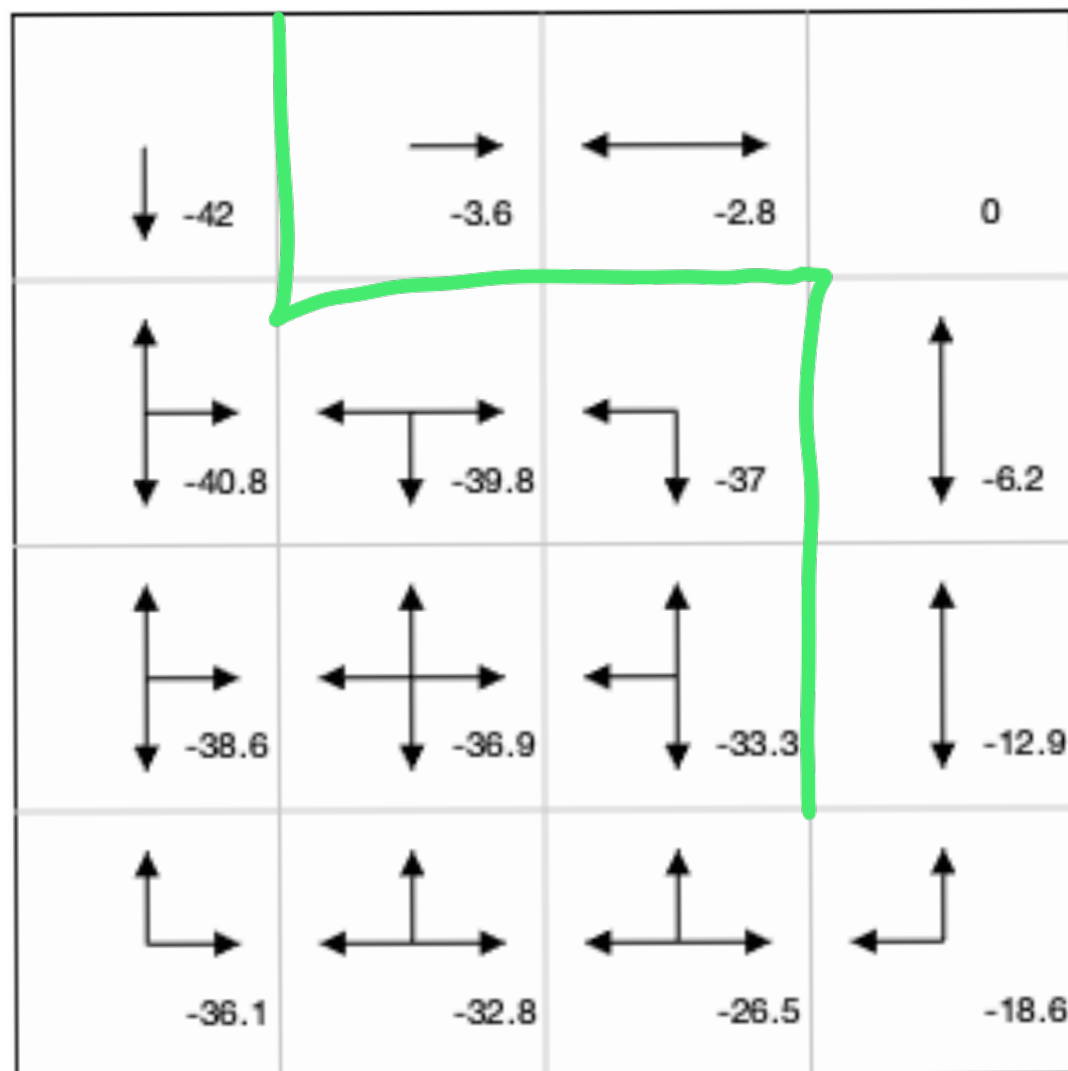
Value functions depend on the policy!
 $V_{\pi}(s)$

Policy Improvement



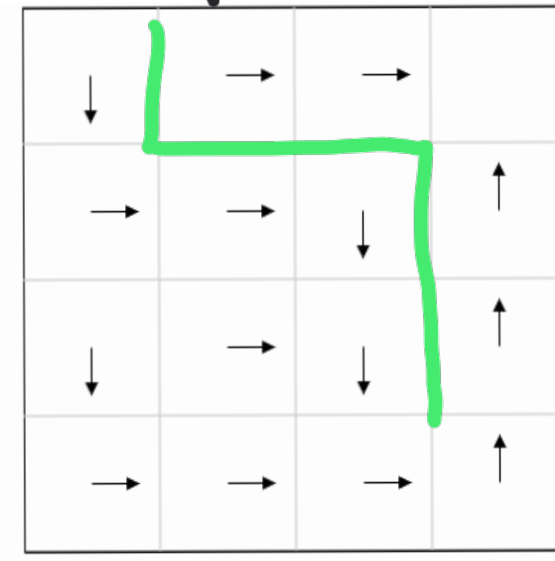
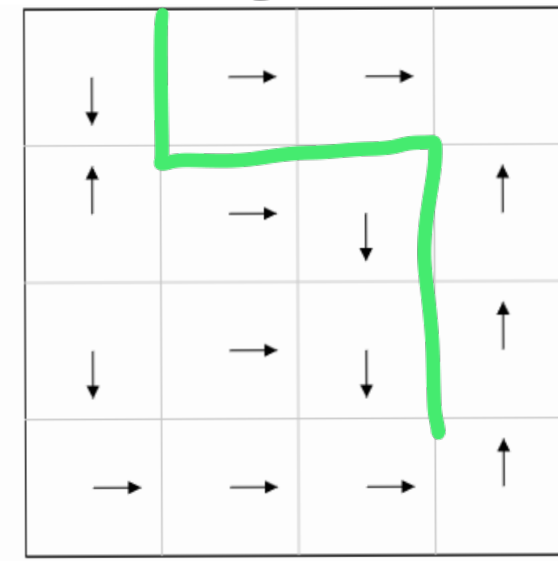
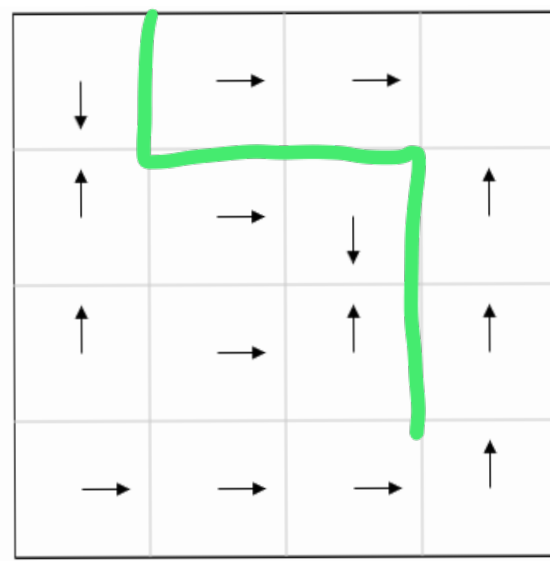
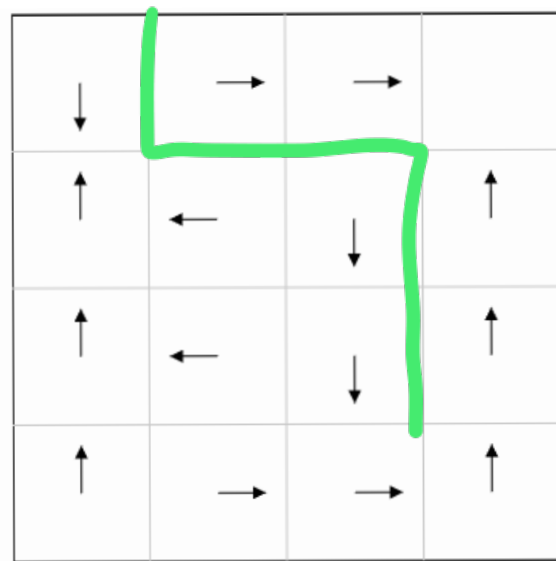
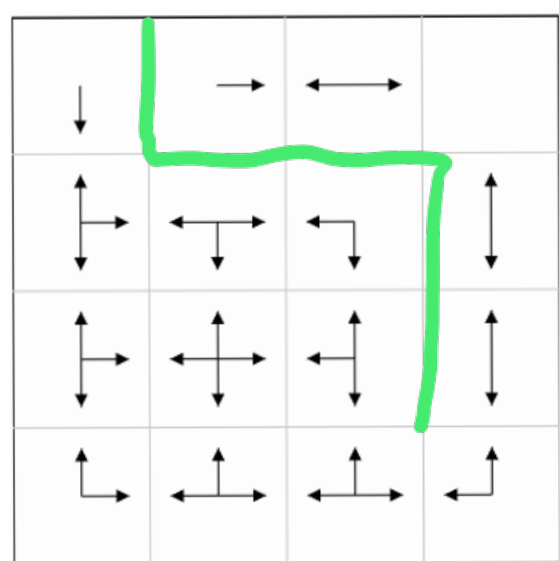
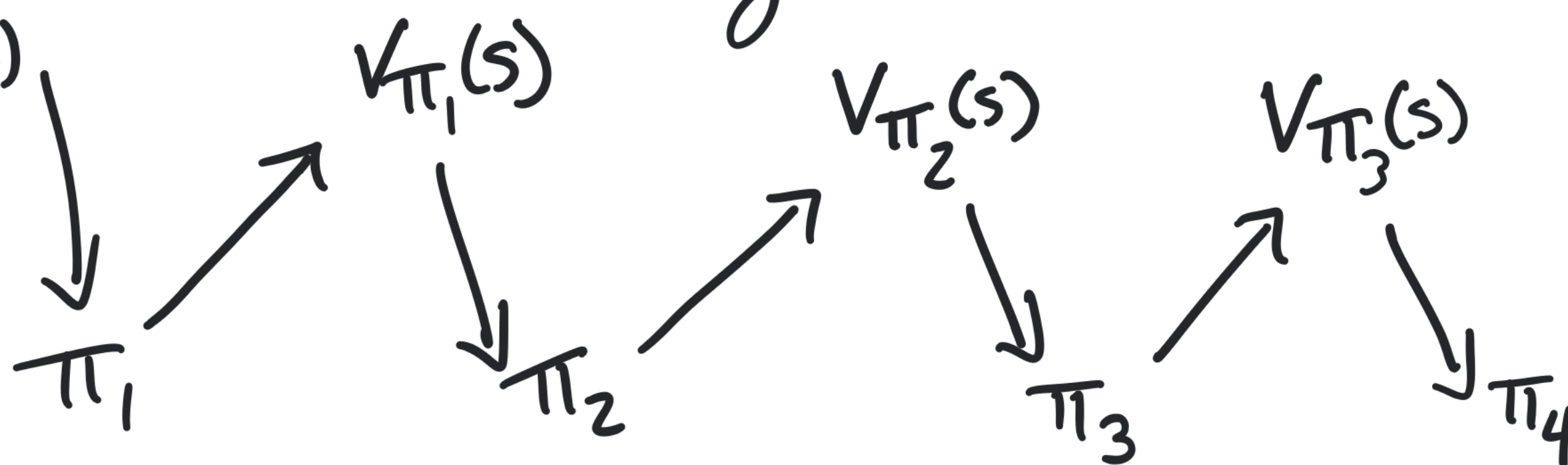
$$\pi_0 = \text{Random}$$

A diagram illustrating a mapping. A curved arrow points from the expression $V_{\pi_0}(s)$ to a box containing a green vertical line.

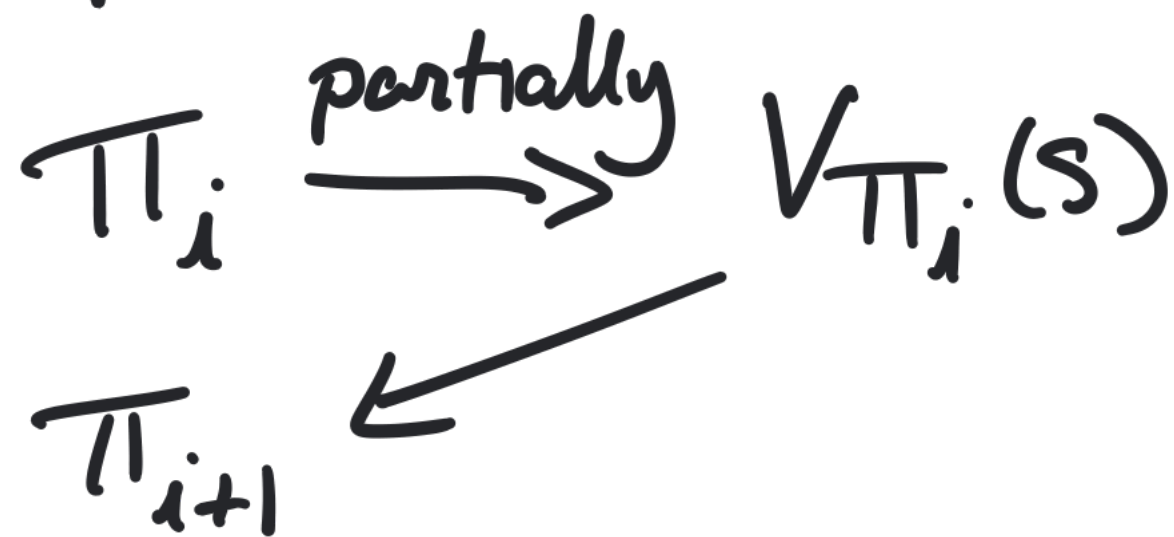
$$\xrightarrow{\quad} \pi_1 \xrightarrow{\quad} V_{\pi_1}(s)$$


GRIDWORLD with A Bad policy

π_0 is $\rightarrow V_{\pi_0}(s)$
biased toward
the bottom
left.



Policy improvement can be incremental



General policy
iteration