# GRIDWORLD &
# The Value function

## BIOE 498/598 PJ

# Markov Decision Process (MDP)

States    S     where we are

actions    a     what we choose (next state)

rewards    r     gain/loss from $(s, a)$

Markov Property : a state tells us everything we need to know (memoryless)

$$s_0, a_0, r_1, s_1, a_1, r_2, s_2, \ldots \ldots, s_T$$

# Making Decisions in an MDP

Assume: MDP is deterministic

Process terminates (@ time $T$)

Let's say we are at state $s_0$.

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

immediate reward

future reward

Value function

# GRIDWORLD

| | | | |
|---|---|---|---|
| 13 | 14 | 15 | finish 16 |
| 9 | 10 | 11 | 12 |
| 5 | 6 | 7 | 8 |
| start 1 | 2 | 3 | 4 |

Compute $V(s)$ for every square

If I know $V(s)$, then I can always find the shortest path.

Get from start to finish in the fewest number of steps.

How do I compute $V(s)$?

$r_i = -1$ except at the finish.

States $S = \{1, 2, \ldots, 16\}$

$\uparrow$ start $\qquad$ $\uparrow$ terminal

actions are defined for every state

| 13 | 14 | 15 | 16 |
|----|----|----|----|
| 9  | 10 | 11 | 12 |
| 5  | 6  | 7  | 8  |
| 1  | 2  | 3  | 4  |

Approach: Monte Carlo Method

1. Pick a state.
2. Take a random walk until I reach the finish
3. Count steps.
4. Repeat.
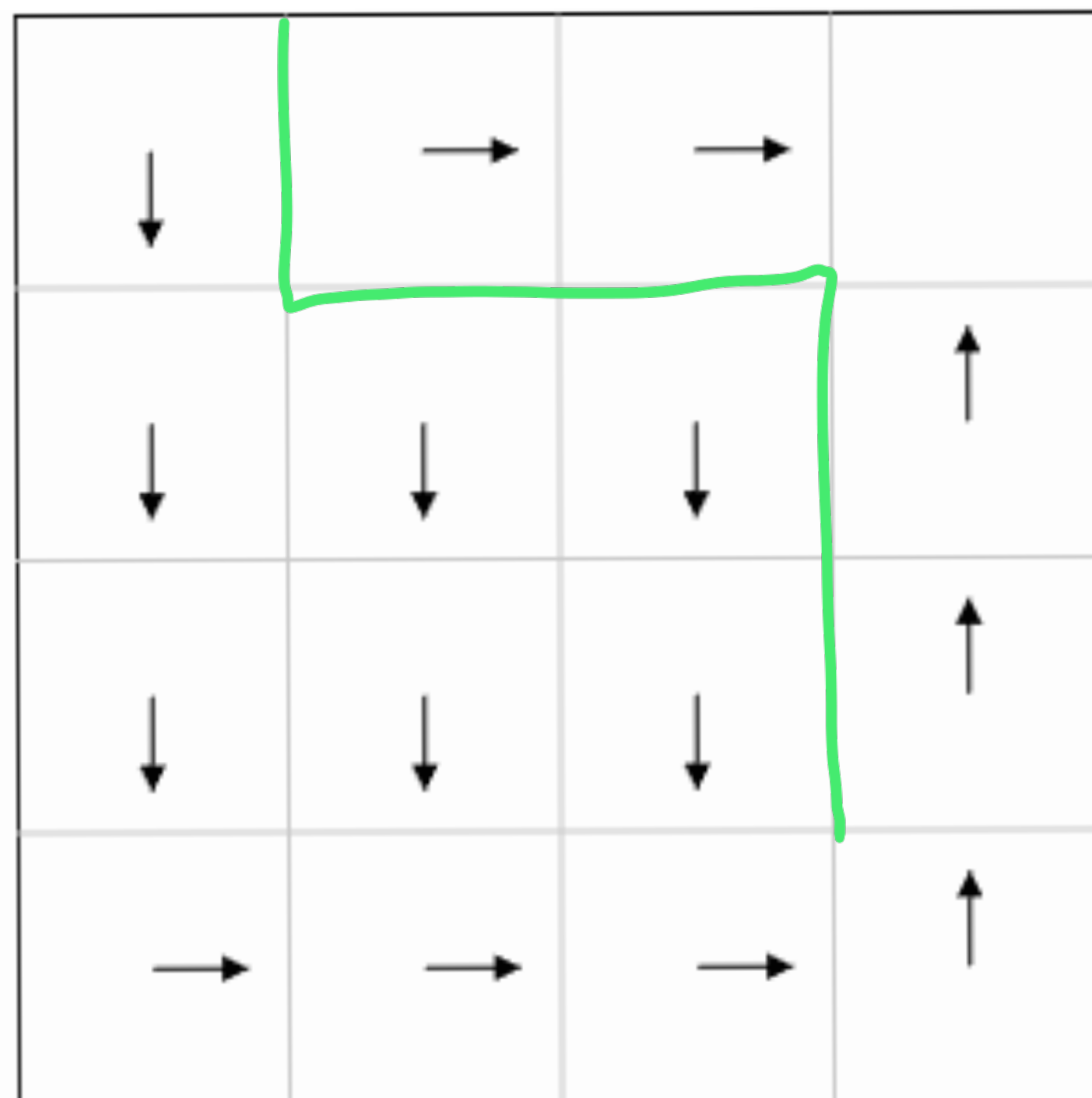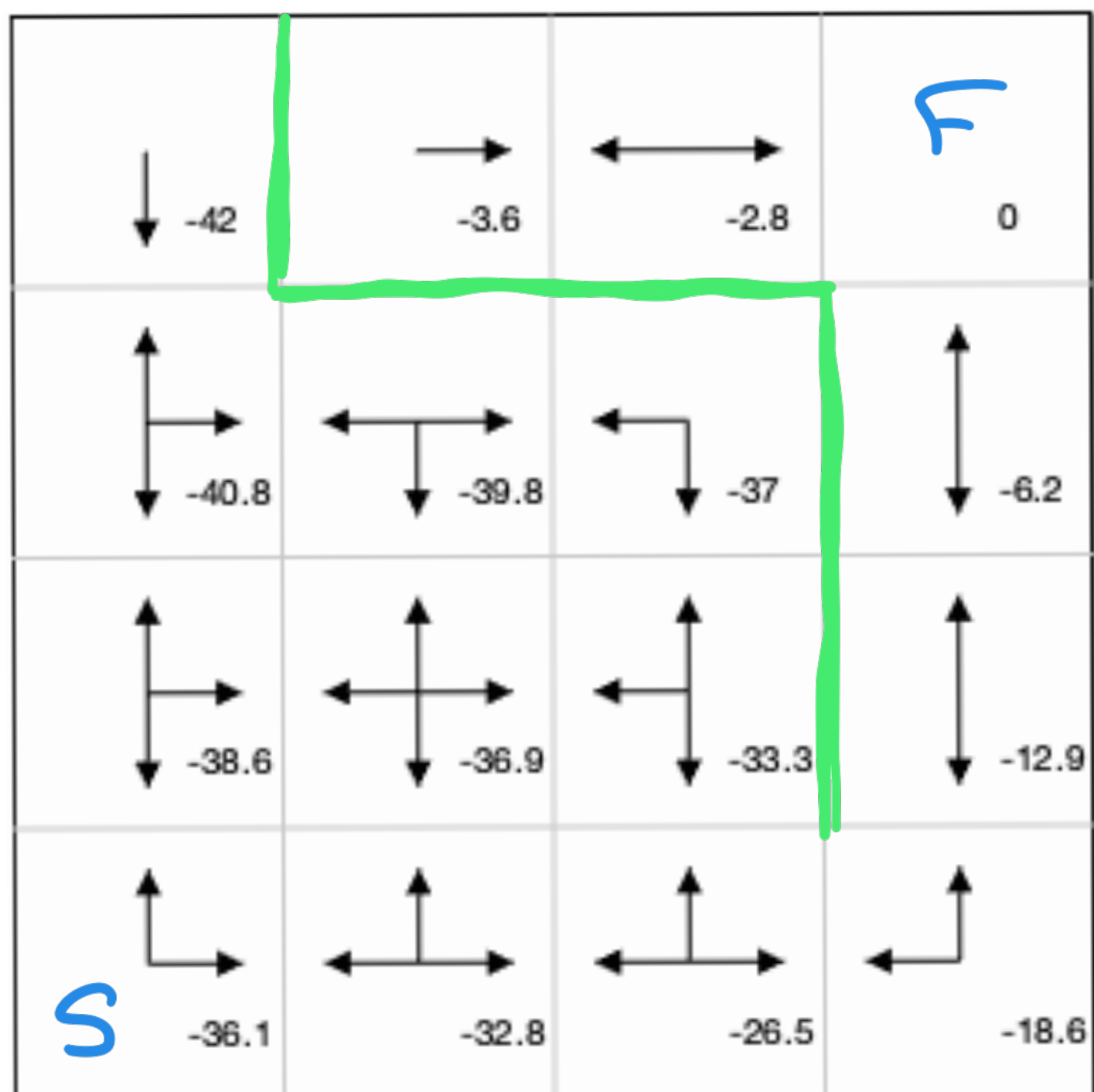
$5, 6, 10, 9, 5, 6, 7, 11, 15, 16$

$V(5) = -9$

# V(s)



| | | | Finish |
|---|---|---|---|
| -25.8 | -22.7 | -14.6 | 0 |
| -27.8 | -26.5 | -21.1 | -14.3 |
| -29.3 | -29.3 | -25.7 | -22.9 |
| -30.2 | -30.4 | -28.3 | -26.5 |

start

# Policy Improvement

V(s)