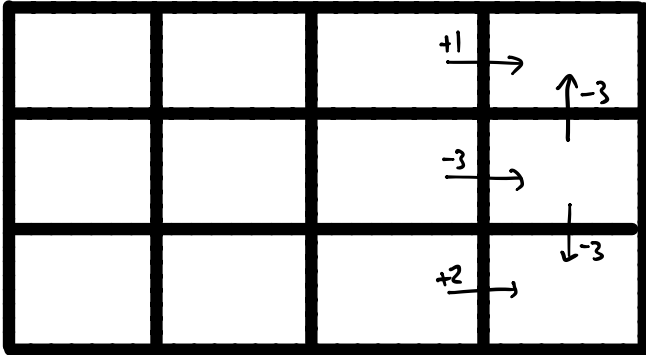
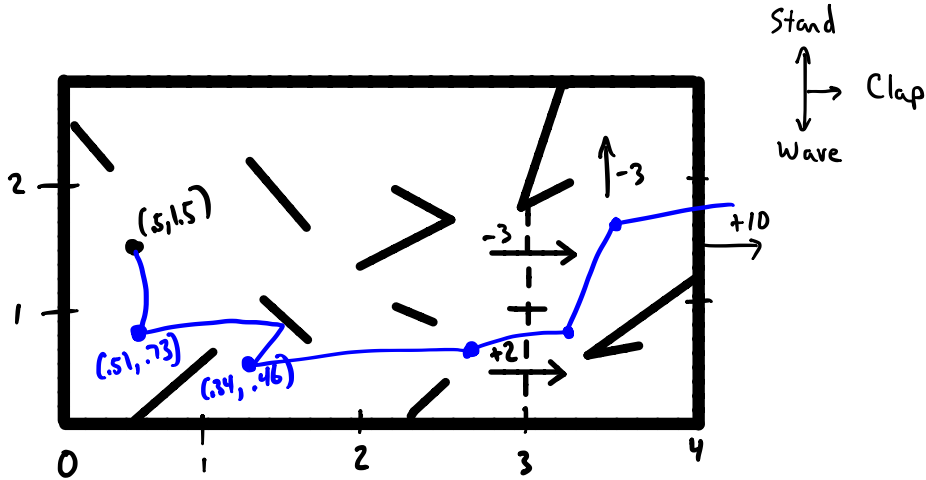
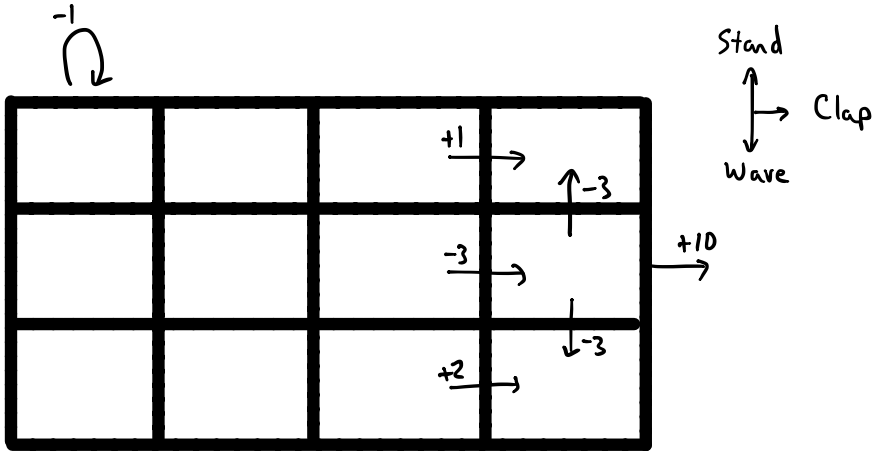


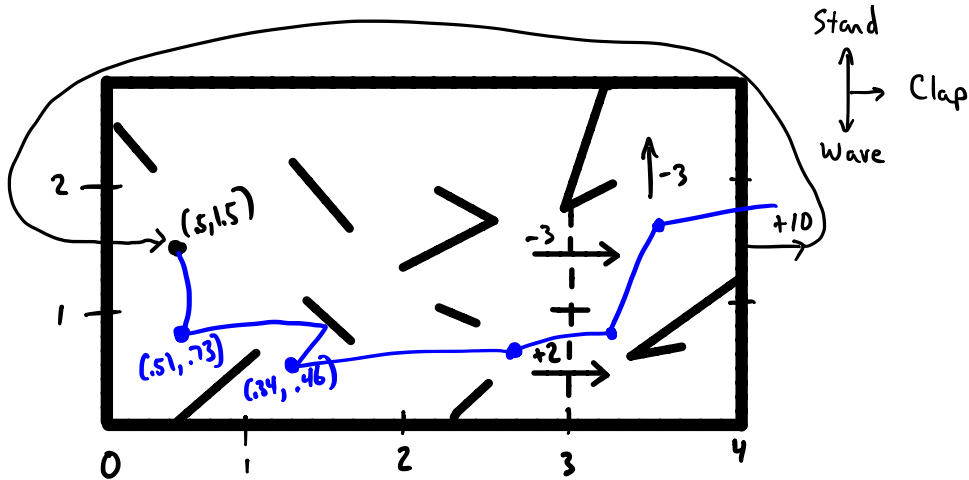
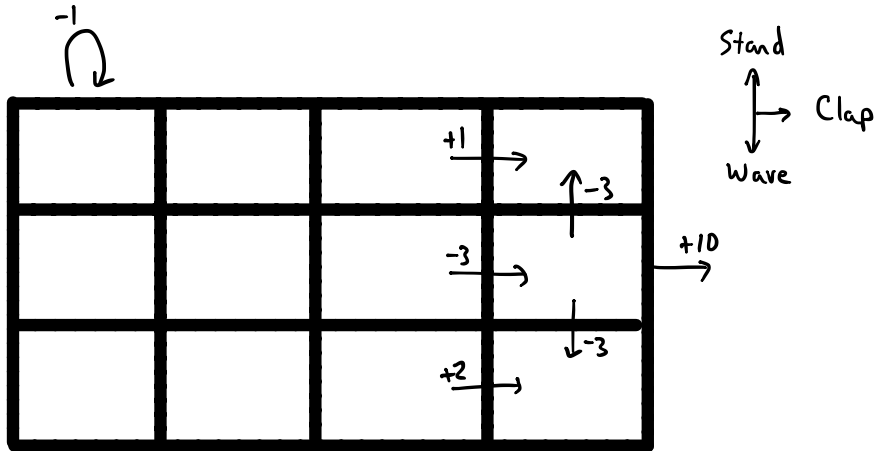
-1



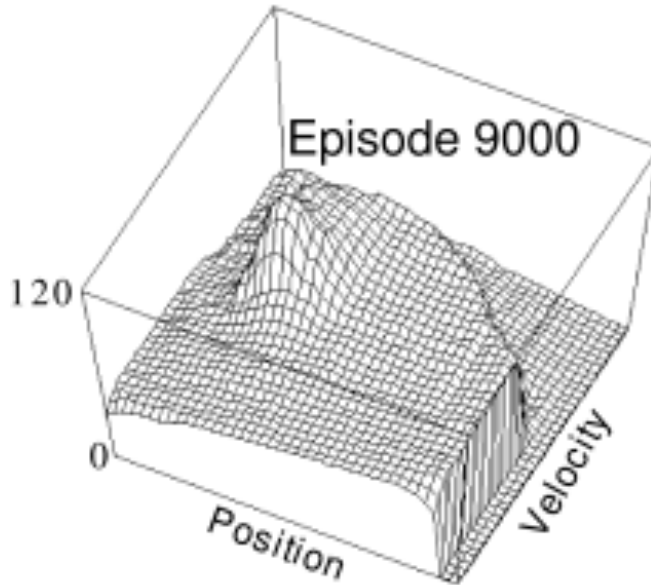
Stand
↑
↓
Wave

Clap
→



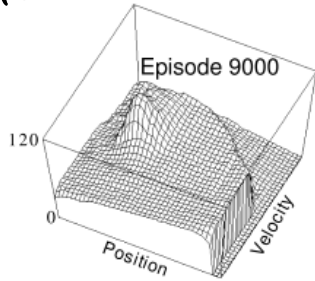


Continuous state: $V(s)$

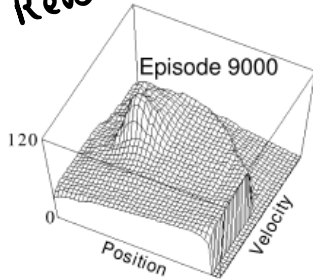


Continuous state: $Q(s,a)$

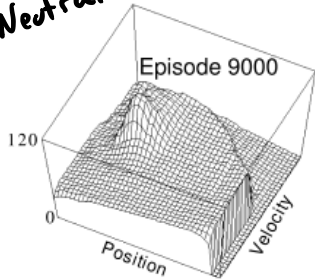
Forward



Reverse



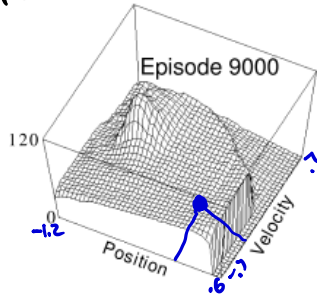
Neutral



Continuous state: $Q(s, a)$

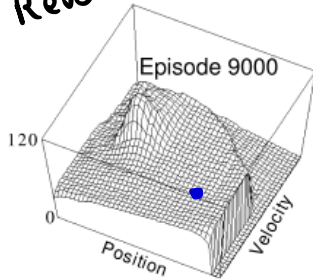
pos = .3, vel = -.3

Forward



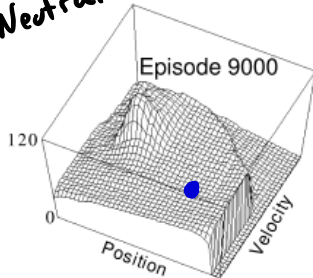
$$Q(s, \text{forward}) = 60$$

Reverse



$$Q(s, \text{reverse}) = 75$$

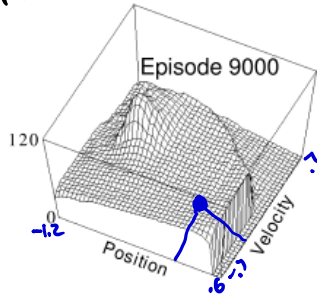
Neutral



$$Q(s, \text{neutral}) = 68$$

Continuous state: $Q(s, a)$; Discrete actions

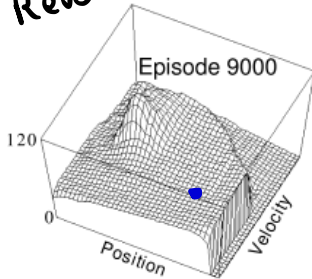
Forward



$$\text{pos} = .3, \text{vel} = -.3$$

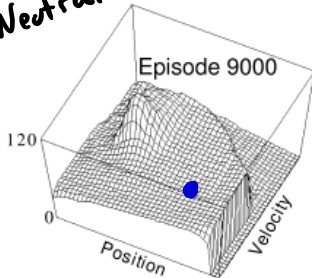
$$Q(s, \text{forward}) = 60 \leftarrow$$

Reverse



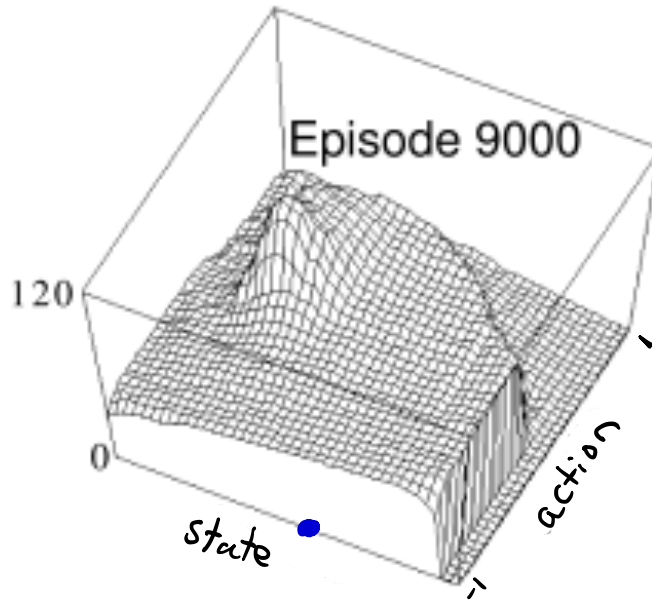
$$Q(s, \text{reverse}) = 75$$

Neutral

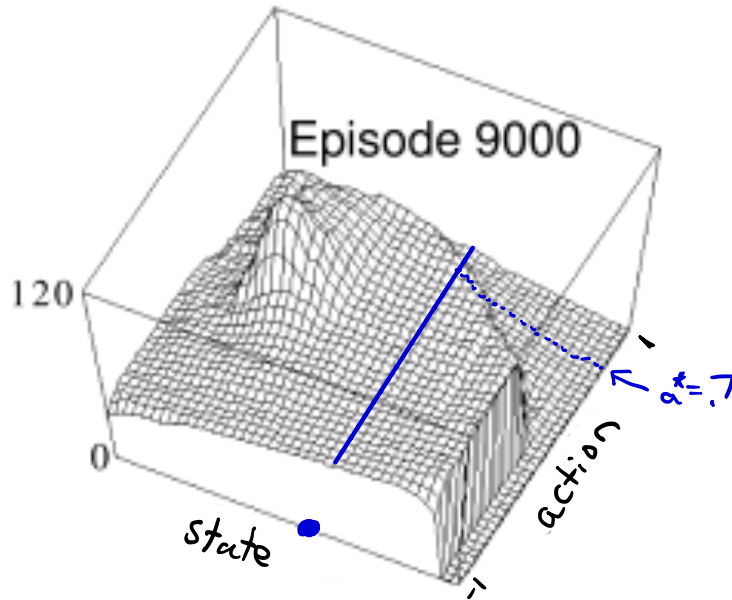


$$Q(s, \text{neutral}) = 68$$

Continuous state, continuous action: $Q(s,a)$



Continuous state, continuous action: $Q(s,a)$

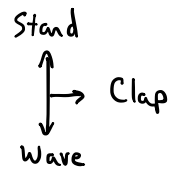
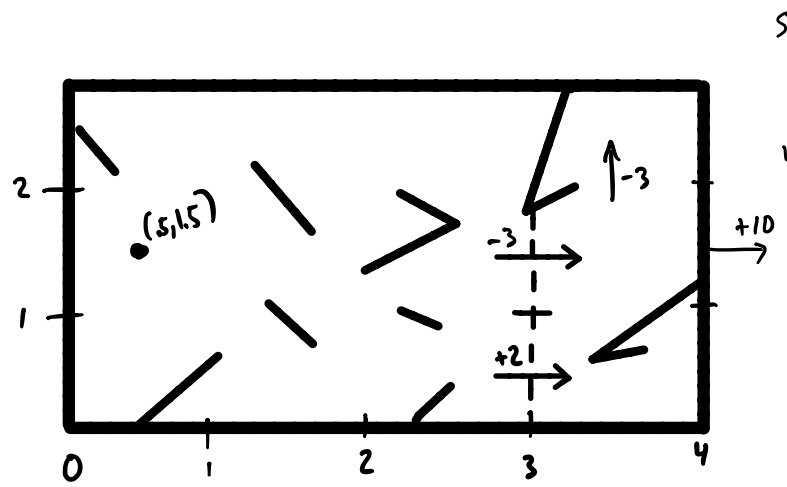
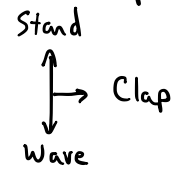
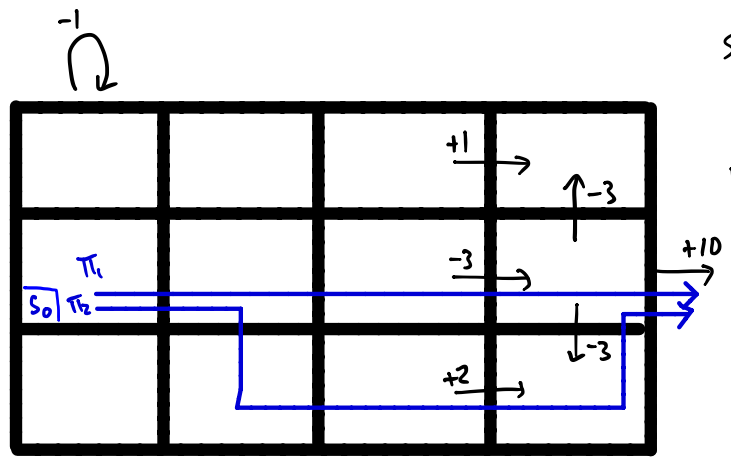


Episodic tasks

Discounting: γ

$$V_{\pi_1}(s_0) =$$

$$V_{\pi_2}(s_0) =$$



Episodic tasks

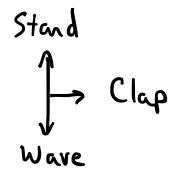
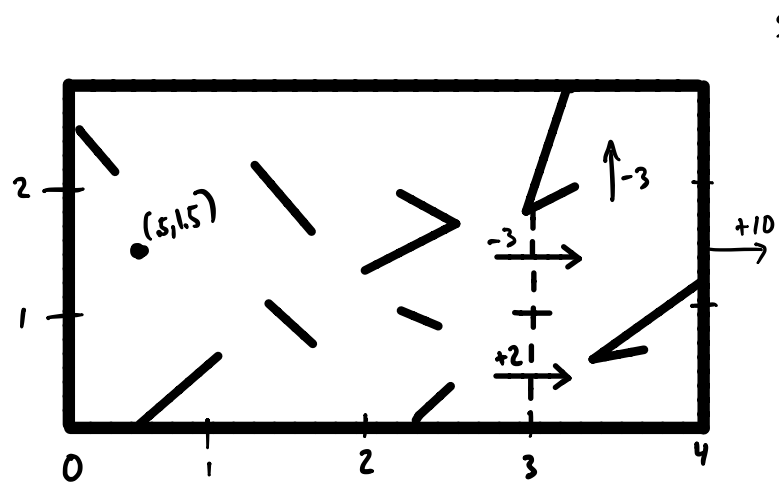
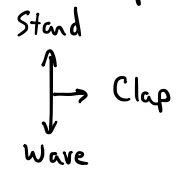
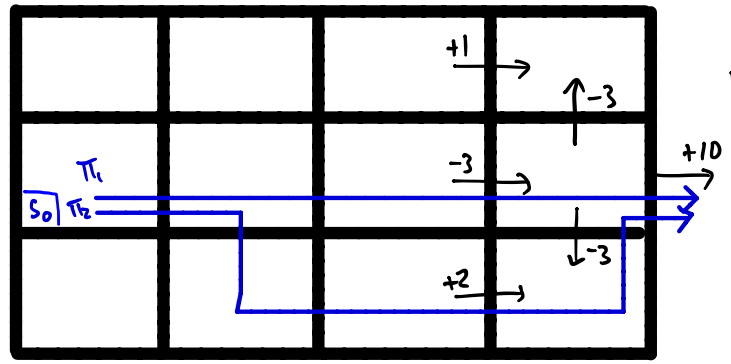
Discounting: γ

$$V_{\pi_1}(s_0) = 0 + 0 - 3\gamma^2 + 10\gamma^3$$

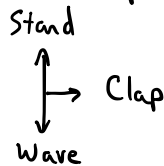
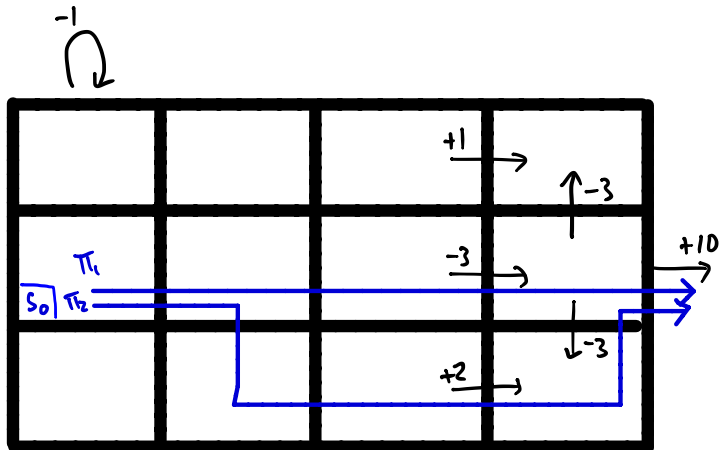
$$V_{\pi_2}(s_0) = 0 + 0 + 0 + 2\gamma^3 + 0 + 10\gamma^5$$

which policy is better?

-1 ↻



Episodic tasks



Discounting: γ

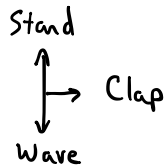
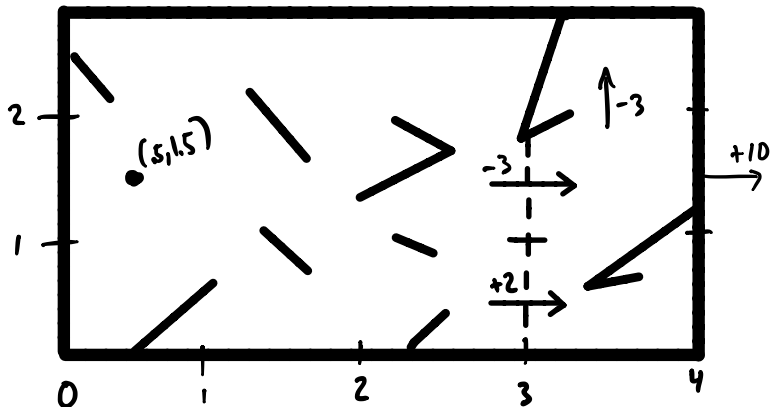
$$V_{\pi_1}(s_0) = 0 + 0 - 3\gamma^2 + 10\gamma^3$$

$$V_{\pi_2}(s_0) = 0 + 0 + 0 + 2\gamma^3 + 0 + 10\gamma^5$$

Which policy is better?

$$\gamma = 1: V_{\pi_1}(s_0) = 7 \quad V_{\pi_2}(s_0) = 12$$

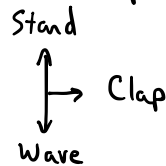
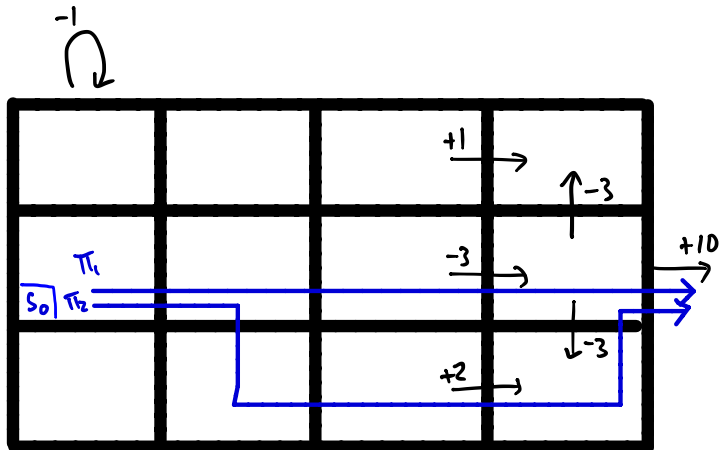
$$\gamma = .5: V_{\pi_1}(s_0) = 1.175 \quad V_{\pi_2}(s_0) = .5625$$



Two meaning of γ :

- 1)
- 2)

Episodic tasks



Discounting: γ

$$V_{\pi_1}(s_0) = 0 + 0 - 3\gamma^2 + 10\gamma^3$$

$$V_{\pi_2}(s_0) = 0 + 0 + 0 + 2\gamma^3 + 0 + 10\gamma^5$$

Which policy is better?

$\gamma = 1$: $V_{\pi_1}(s_0) = 7$ $V_{\pi_2}(s_0) = 12$

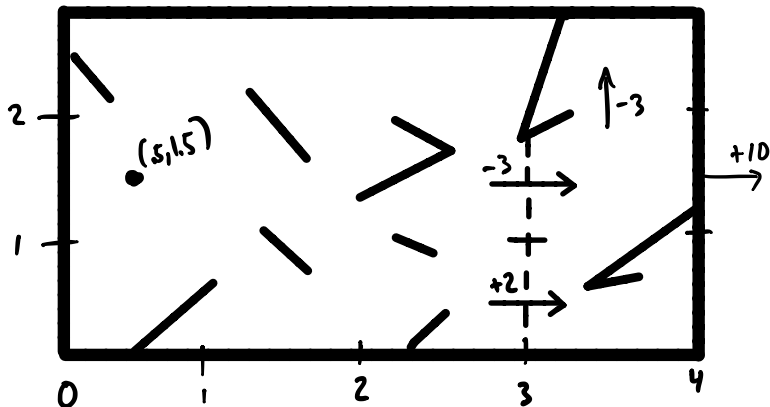


$\gamma = .5$: $V_{\pi_1}(s_0) = 1.175$ $V_{\pi_2}(s_0) = .5625$

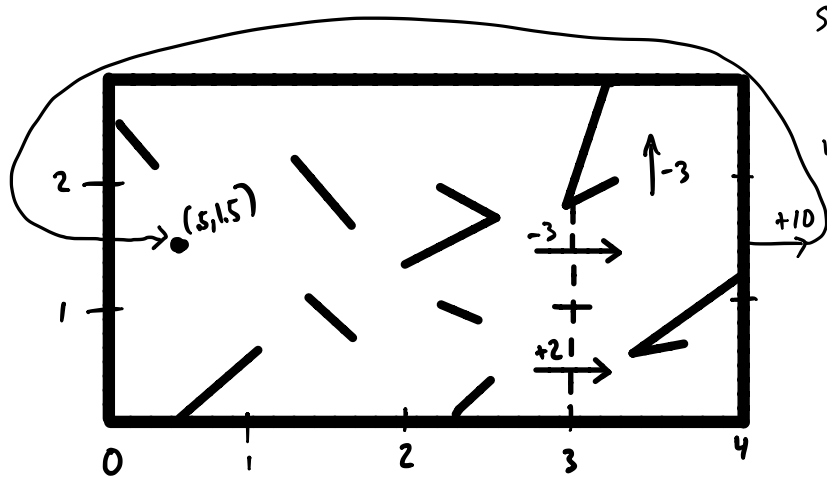
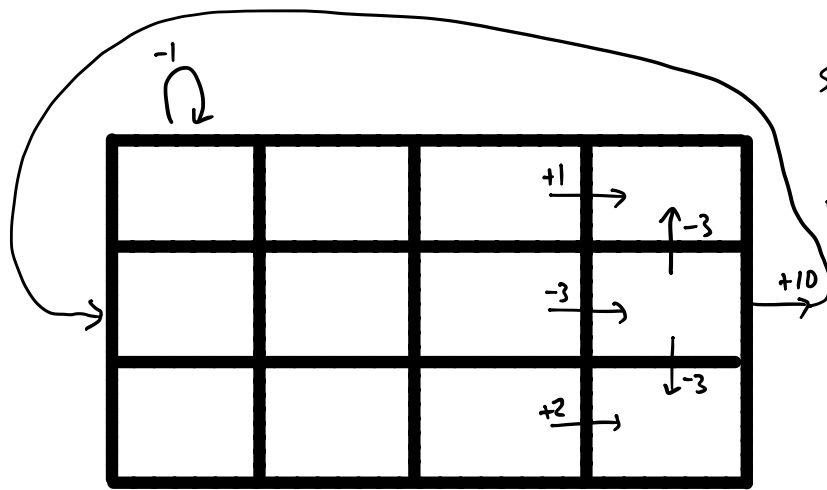
Two meaning of γ :

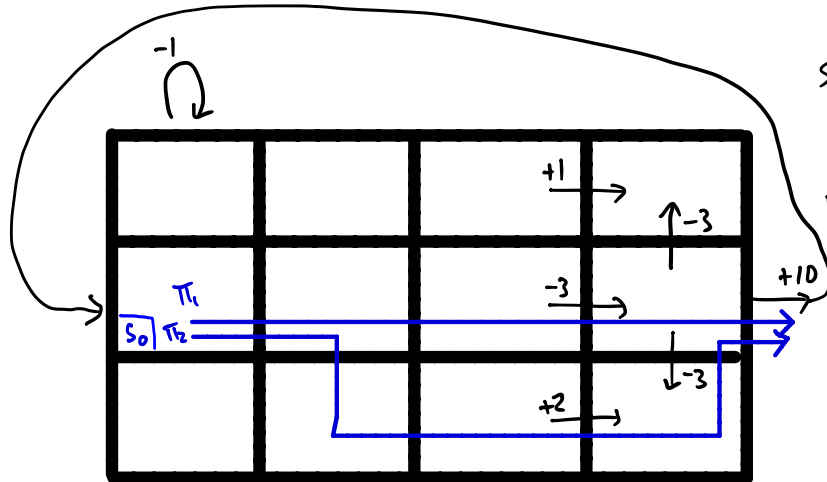
1) interest / inflation

2) probability of episode ending $(1-\gamma)$



Continuing tasks





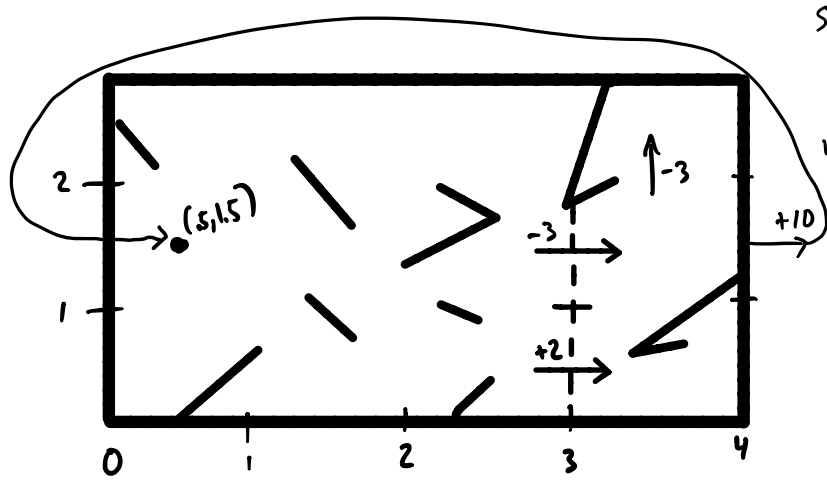
Continuing tasks

Stand \updownarrow Clap
Wave \rightarrow

Discounting: γ

$V_{\pi_1}(s_0) =$

$V_{\pi_2}(s_0) =$



Stand \updownarrow Clap
Wave \rightarrow

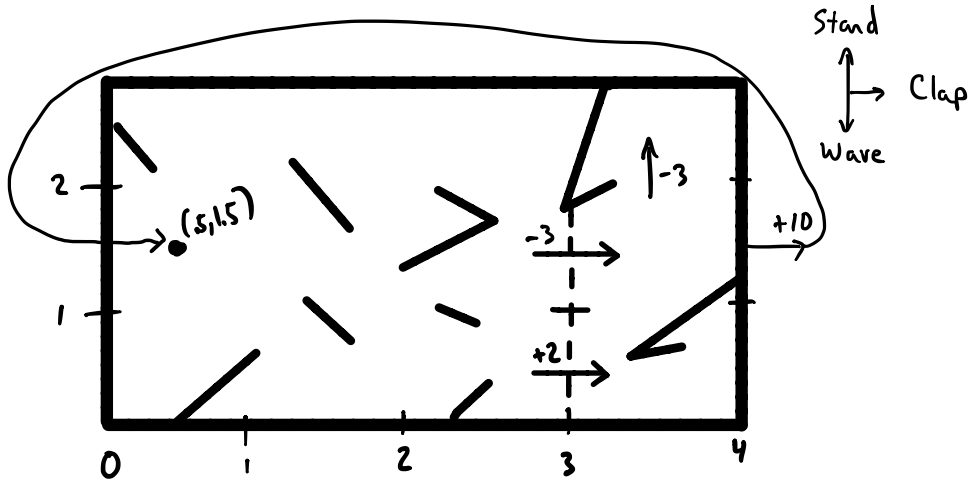
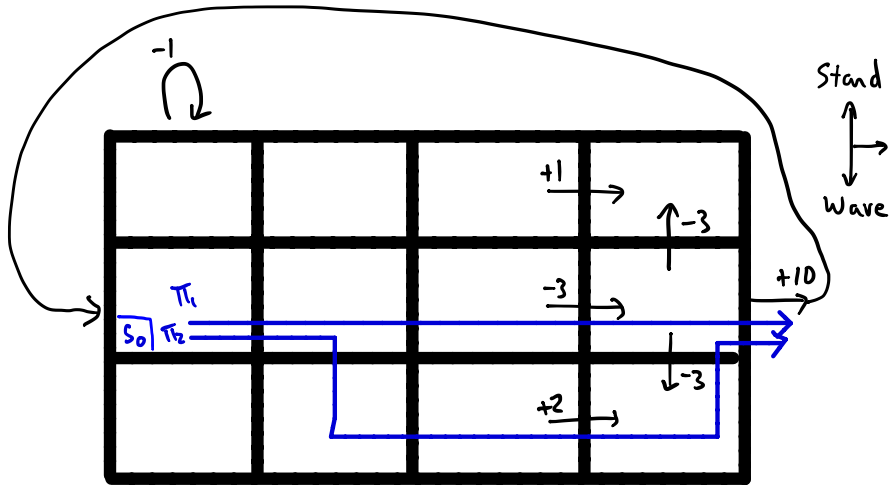
Continuing tasks

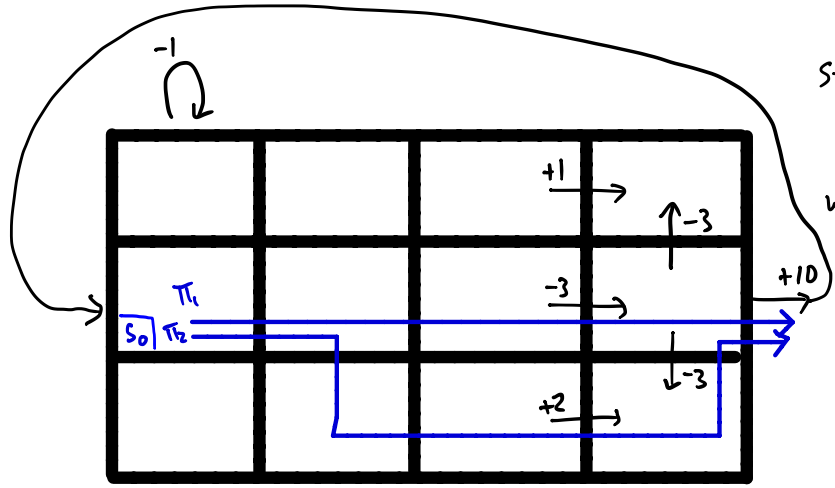
Discounting: δ

$$V_{\pi_1}(s_0) = 0 + 0 - 3\delta^2 + 10\delta^3 + 0 + 0 - 3\delta^6 + 10\delta^7 + 0 + 0 - 3\delta^{10} + 10\delta^{11} + \dots$$

$$V_{\pi_2}(s_0) = 0 + 0 + 0 + 2\delta^3 + 0 + 10\delta^5 + 0 + 0 + 0 + 2\delta^9 + 0 + 10\delta^{11} + \dots$$

which policy is better?





Continuing tasks

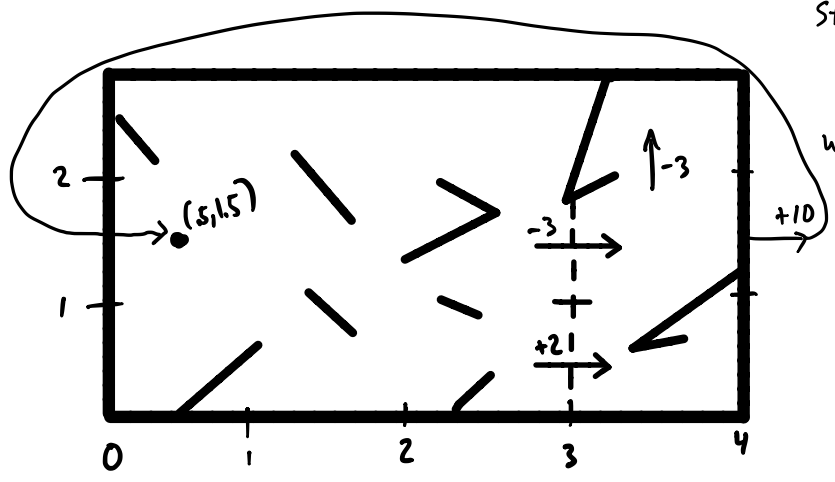
Stand \updownarrow Clap \rightarrow
Wave

Discounting: δ

$$V_{\pi_1}(s_0) = 0 + 0 - 3\delta^2 + 10\delta^3 + 0 + 0 - 3\delta^6 + 10\delta^7 + 0 + 0 - 3\delta^{10} + 10\delta^{11} + \dots$$

$$V_{\pi_2}(s_0) = 0 + 0 + 0 + 2\delta^3 + 0 + 10\delta^5 + 0 + 0 + 0 + 2\delta^9 + 0 + 10\delta^{11} + \dots$$

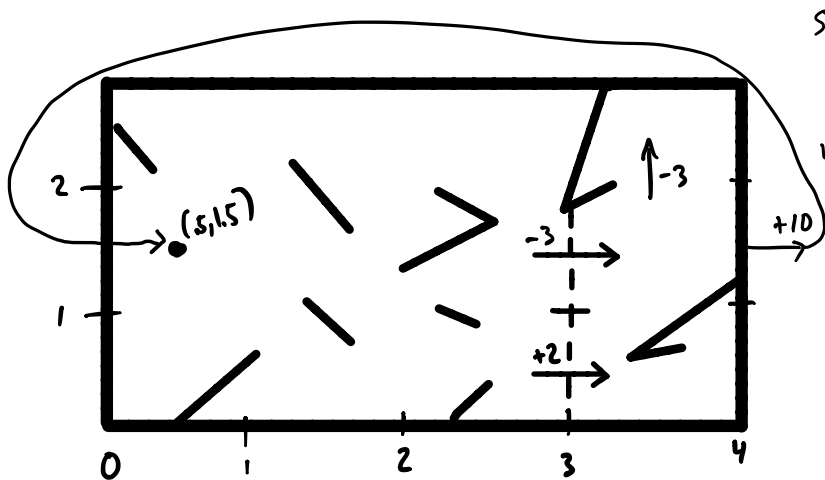
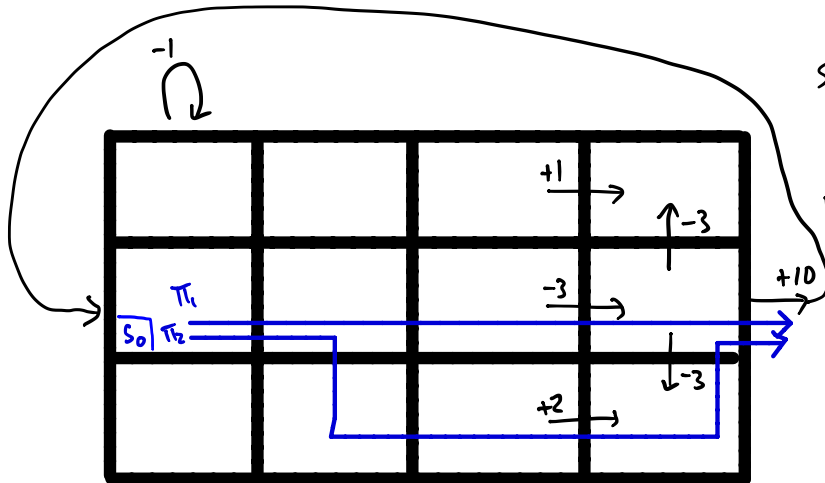
which policy is better?

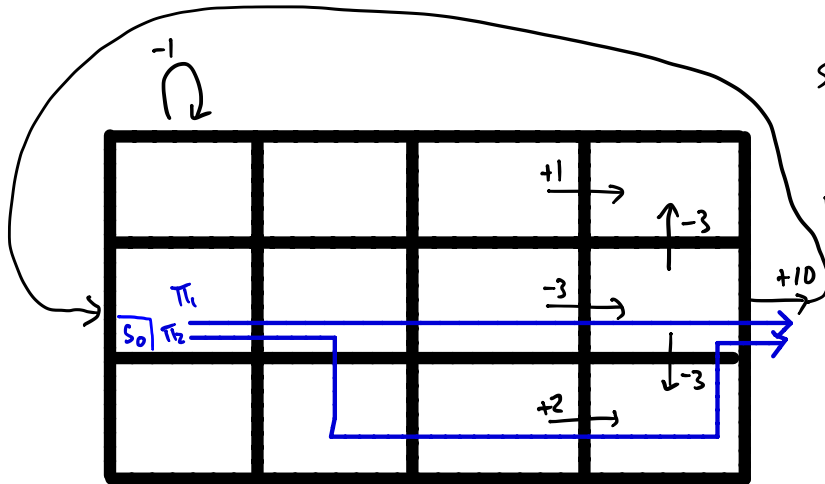


Stand \updownarrow Clap \rightarrow
Wave

Discrete state (tabular):
Depends on δ

Continuous state (function approx.):
Might not depend on δ !





Continuing tasks

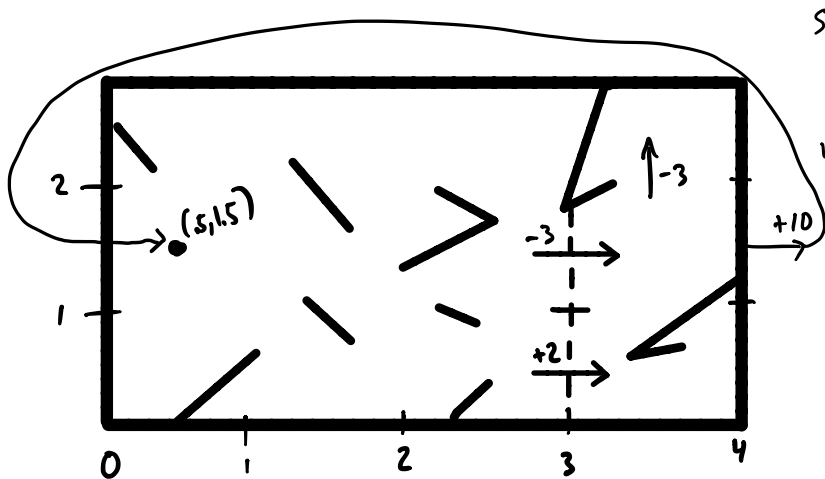
Stand
↕
Clap
↕
Wave

Average reward RL

$$r(\pi_1) = 7/4$$

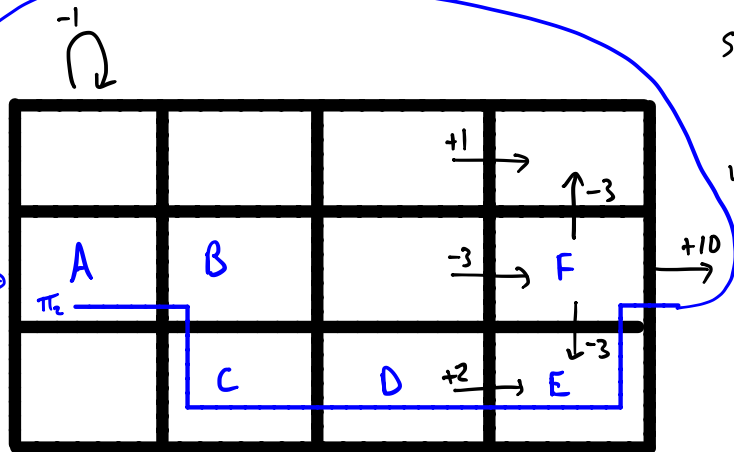
$$r(\pi_2) = 12/6 = 2$$

which policy is better?



Stand
↕
Clap
↕
Wave

Differential value function

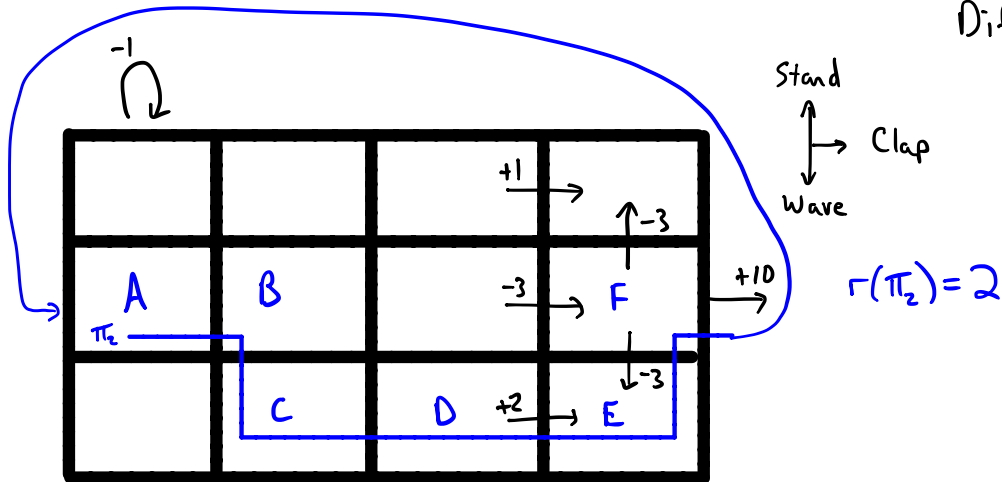


$r(\pi_2) = 2$ ← estimated by algorithm in the book: β
 Differential semi-gradient SARSA (R-learning)

- $V(A) =$
- $V(B) = ?$
- $V(C) = 0$
- $V(D) =$
- $V(E) =$
- $V(F) =$

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{r, s'} p(s', r | s, a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



$$V(A) = ?$$

$$V(B) = -2$$

$$V(C) = 0$$

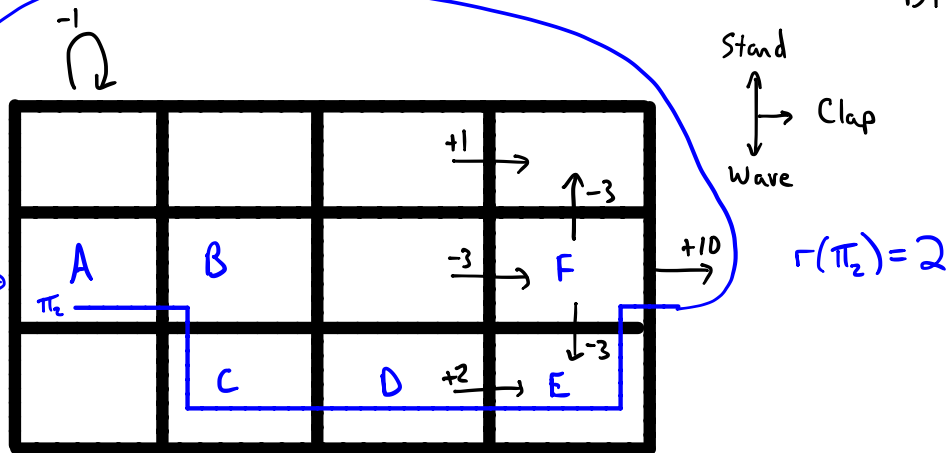
$$V(D) = ?$$

$$V(E) = ?$$

$$V(F) = ?$$

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

Differential value function

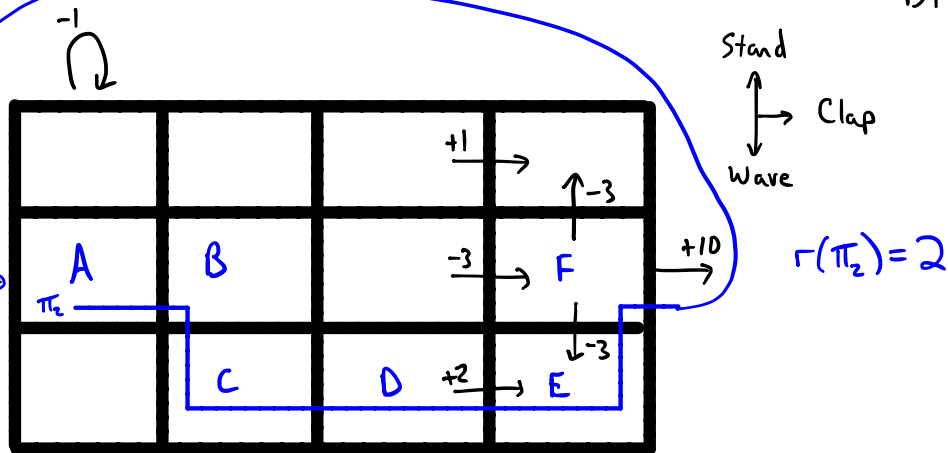


- $V(A) = -4$
- $V(B) = -2$
- $V(C) = 0$
- $V(D) = 2$
- $V(E) = 2$
- $V(F) = 4$

Can this be V?

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{r, s'} \rho(s', r | s, a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



$$\begin{aligned} V(A) &= -4 \\ V(B) &= -2 \\ V(C) &= 0 \\ V(D) &= 2 \\ V(E) &= 2 \\ V(F) &= 4 \end{aligned}$$

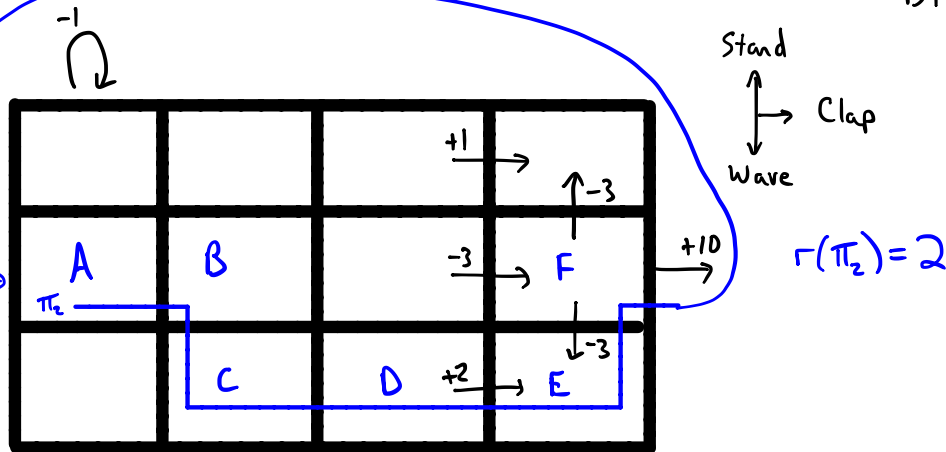
Can this be V ?

$$V(A) + V(B) + \dots + V(F) = 2$$

But avg. value of a cycle
must be 0....

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{r,s'} p(s',r|s,a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



What's the steady state distribution of π_2 ?

$$\begin{aligned} V(A) &= -4 - 1/3 = -13/3 \\ V(B) &= -2 - 1/3 = -7/3 \\ V(C) &= 0 - 1/3 = -1/3 \\ V(D) &= 2 - 1/3 = 5/3 \\ V(E) &= 2 - 1/3 = 5/3 \\ V(F) &= 4 - 1/3 = 11/3 \end{aligned}$$

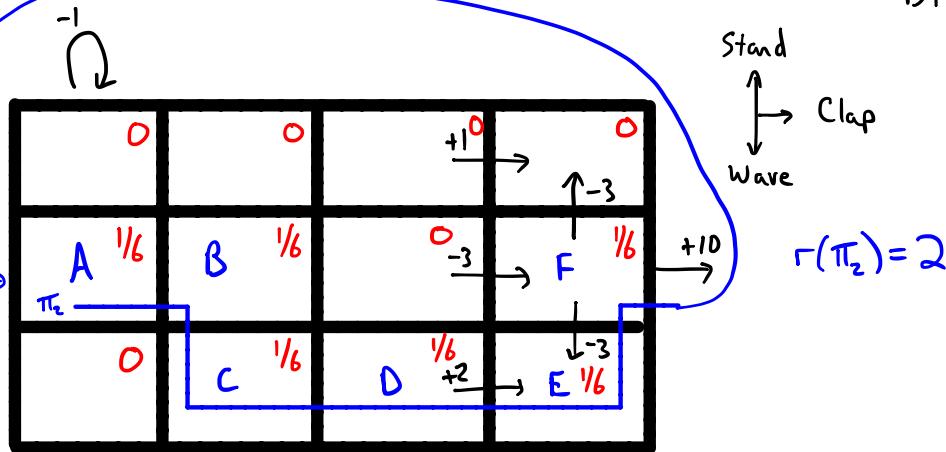
Can this be V ?

$$V(A) + V(B) + \dots + V(F) = 2$$

But avg. value of a cycle must be 0....

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} p(s', r|s, a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



$$\begin{aligned} V(A) &= -4 - 1/3 = -13/3 \\ V(B) &= -2 - 1/3 = -7/3 \\ V(C) &= 0 - 1/3 = -1/3 \\ V(D) &= 2 - 1/3 = 5/3 \\ V(E) &= 2 - 1/3 = 5/3 \\ V(F) &= 4 - 1/3 = 11/3 \end{aligned}$$

Can this be V ?

$$V(A) + V(B) + \dots + V(F) = 2$$

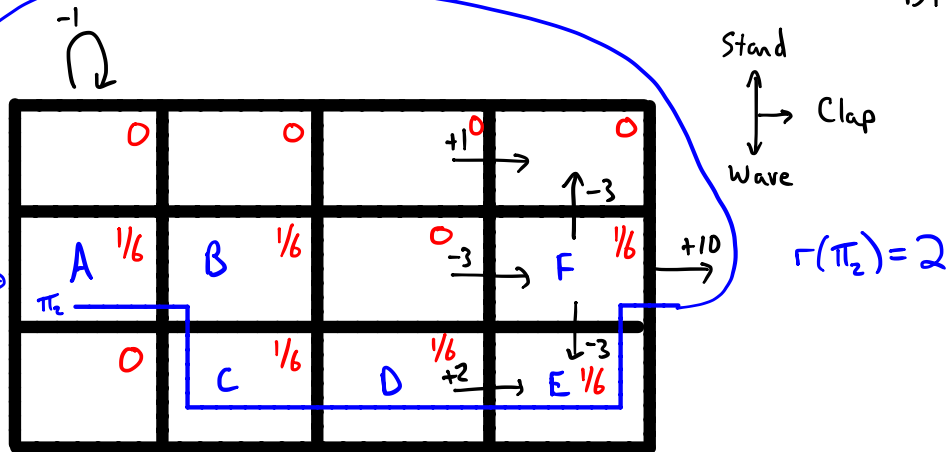
But avg. value of a cycle must be 0....

What's the steady state distribution of π_2 ?

Is this MDP ergodic?
(what does ergodic mean?)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{r,s'} p(s',r|s,a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



$$\begin{aligned} V(A) &= -4 - 1/3 = -13/3 \\ V(B) &= -2 - 1/3 = -7/3 \\ V(C) &= 0 - 1/3 = -1/3 \\ V(D) &= 2 - 1/3 = 5/3 \\ V(E) &= 2 - 1/3 = 5/3 \\ V(F) &= 4 - 1/3 = 11/3 \end{aligned}$$

Can this be V ?

$$V(A) + V(B) + \dots + V(F) = 2$$

But avg. value of a cycle must be 0....

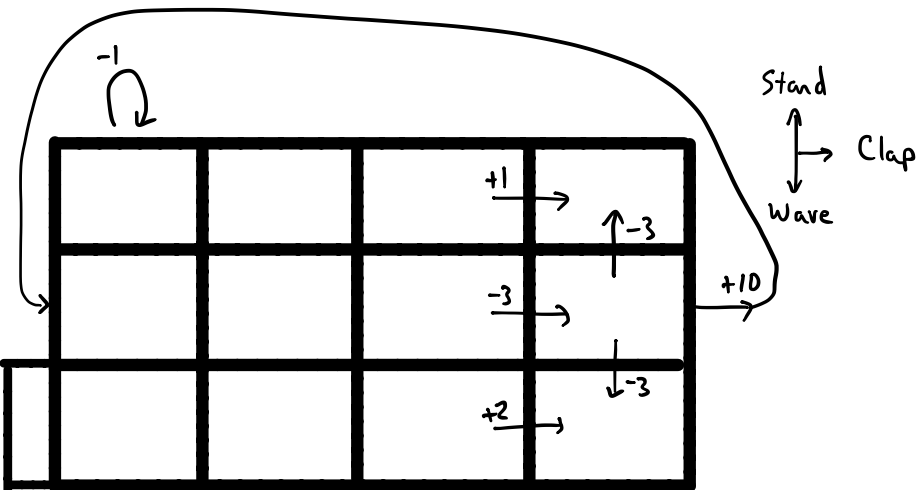
What's the steady state distribution of π_2 ?

Is this MDP ergodic?

(i.e. does every policy have a steady state distribution independent of S_0 ?)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



What about this MDP?

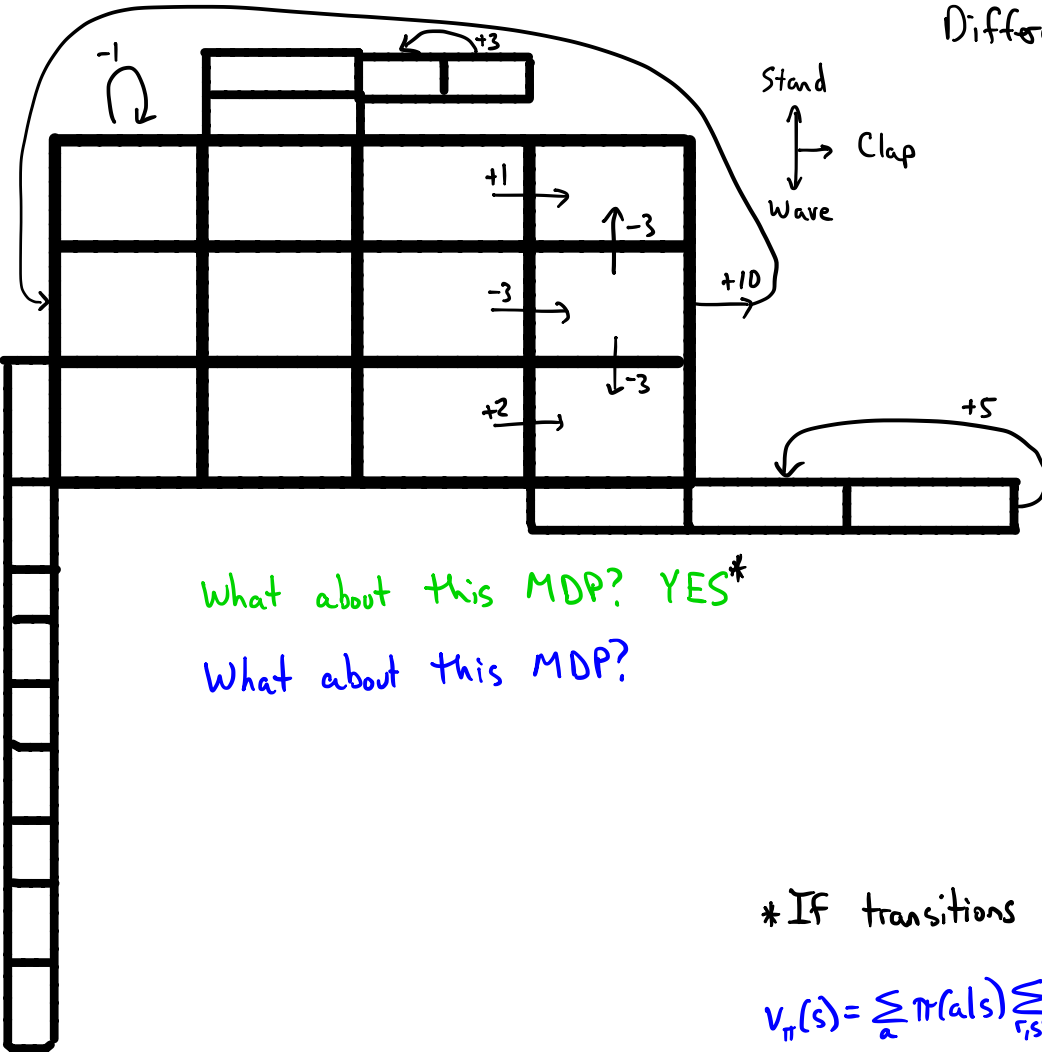
What's the steady state distribution of π_2 ?

Is this MDP ergodic? YES*

* If transitions are at least slightly stochastic (\updownarrow)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



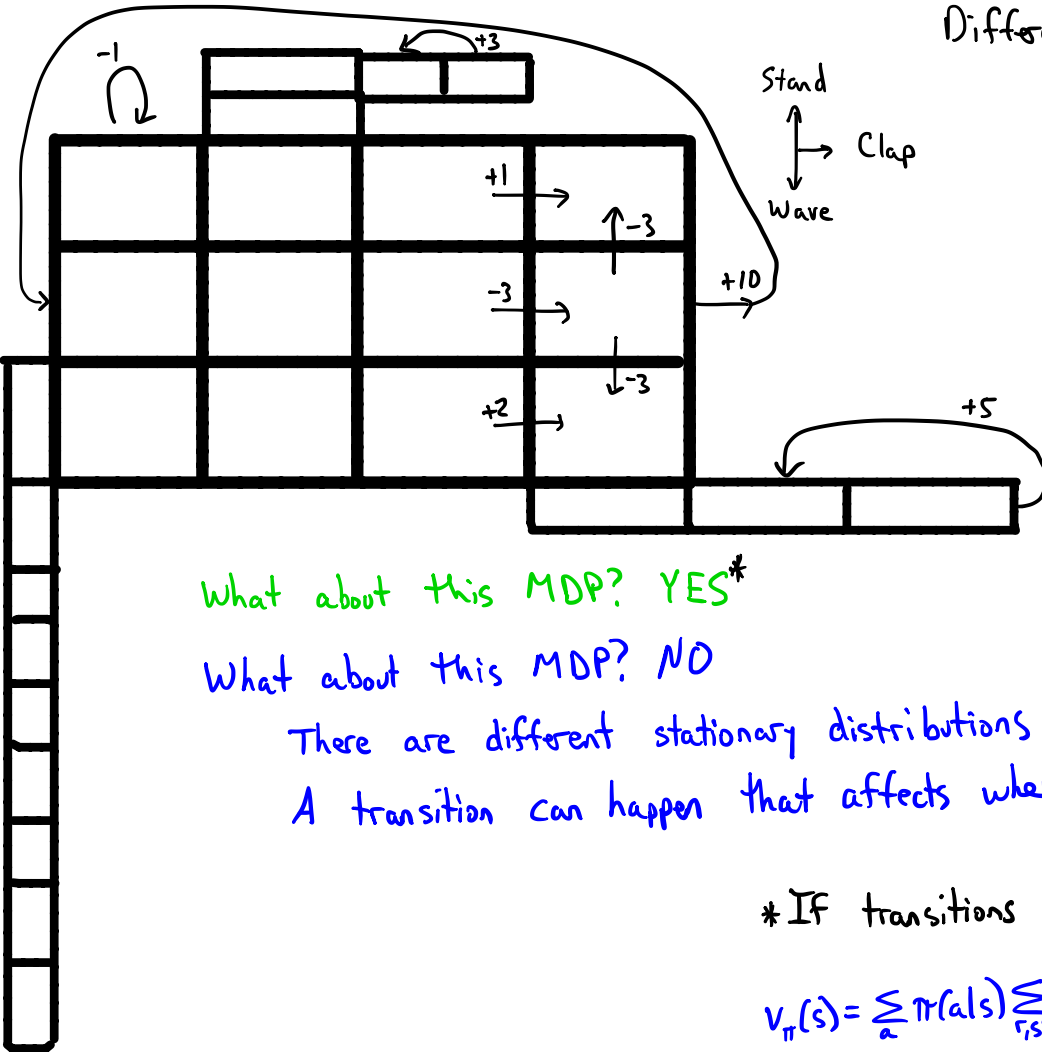
What's the steady state distribution of π_2 ?
 Is this MDP ergodic? YES

What about this MDP? YES*
 What about this MDP?

* IF transitions are at least slightly stochastic (↕)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r - r(\pi) + V_{\pi}(s')]$$

Differential value function



What's the steady state distribution of π_2 ?
 Is this MDP ergodic? YES

What about this MDP? YES*

What about this MDP? NO

There are different stationary distributions for the same policy
 A transition can happen that affects where the agent can get (eventually)

* IF transitions are at least slightly stochastic (\updownarrow)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{r, s'} p(s', r | s, a) [r - r(\pi) + V_{\pi}(s')]$$