Stand

Clap

Wave

Stand

Clap

Wave

-1

| $0$ | $0$ | $0$ | $+1 \rightarrow$ $0$ |
|---|---|---|---|
| $0$ | $0$ | $0$ | $-3 \rightarrow$ $0$ |
| $0$ | $0$ | $0$ | $+2 \rightarrow$ $0$ |

$-3$ (↑)

$+10 \rightarrow$

$-3$ (↓)

Stand

Clap

Wave

Stand

Clap

Wave

-1

Stand
Clap
Wave

| 0 | 0 | 0 | +1 → | 0 |
| 0 | 0 | 0 | -3 → | 0 |
| 0 | 0 | 0 | +2 → | 0 |

-3
+10
-3

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

| | | .5 | -.25 |
|---|---|---|---|
| O | O | 1 | 7.5 |
| O | O | O | O |

| | | | |
|---|---|---|---|
| | | | |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

$S, -1$

$C, +1$

$C, 0$

$W, 0$

$C, -3$

$C, +10$

| 0 | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

| | | | |
|---|---|---|---|
| | | | ? |
| | | | |

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

$S, -1$

$C, +1$

$W, 0$

$C, 0$

$C, -3$

$C, +10$

| | | | |
|---|---|---|---|
| O | O | .5 | -.25 |
| O | O | 1 | 7.5 |
| O | O | O | O |

| | | | |
|---|---|---|---|
| | | ? | ? |
| | | ? | ~10 |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

$\alpha = .5, \ \gamma = 1$

$S, -1$

Stand

$C, +1$

Clap

$W, 0$

Wave

$C, 0$

$C, -3$

$C, +10$

| 0 | 0 | .5 | -.25 |
|---|---|-----|------|
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

← on policy *

| | | | ~9.5 |
|---|---|---|------|
| | | | ~10 |
| | | | |

* Dyna-Q learns Q values
so $(\uparrow, 9)$ and $(\downarrow, 10)$

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

$S,-1$  $C,+1$  $W,0$  $C,-3$  $C,+10$  $C,0$

| 0 | 0 | .5 | -.25 |
|---|---|----|------|
| 0 | 0 | 1  | 7.5  |
| 0 | 0 | 0  | 0    |

on policy

| | | ~10.5 | ~9.5 |
|---|---|---|---|
| ~8.75 | ~8.75 | ~8.75 | ~10 |
| | | | |

| | | | |
|---|---|---|---|
| | | | |
| | | | |

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

$S, -1$

$C, +1$

$W, 0$

$C, +10$

$C, 0$

$C, -3$

on policy

| 0 | 0 | .5 | -.25 |
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

| | | ~10.5 | ~9.5 |
| ~8.75 | ~8.75 | ~8.75 | ~10 |
| | | | |

| | | | |
| | | | |
| | | | |

— What about unvisited states?

— What if transition function were stochastic?

$\alpha = .5, \, \gamma = 1$

Stand

Clap

Wave

$s, -1$

$c, +1$

$w, 0$

$c, 0$

$c, -3$

$c, +10$

$c, +2$

$s, 0$

| | | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

← on policy

| 0 | 0 | ~10.5 | ~9.5 |
|---|---|---|---|
| ~8.75 | ~8.75 | ~8.75 | ~10 |
| 0 | 0 | 0 | 0 |

| | | | |
|---|---|---|---|
| | | | |
| | | | |

— What about unvisited states?

— What if transition function were stochastic?

— Does the order of updates matter?

$\alpha = .5, \ \gamma = 1$

Stand

Clap

Wave

$S, -1$

$C, +1$

$W, 0$

$C, +10$

$C, 0$

$C, -3$

$C, +2$

$S, 0$

on policy

| 0 | 0 | .5 | -.25 |
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

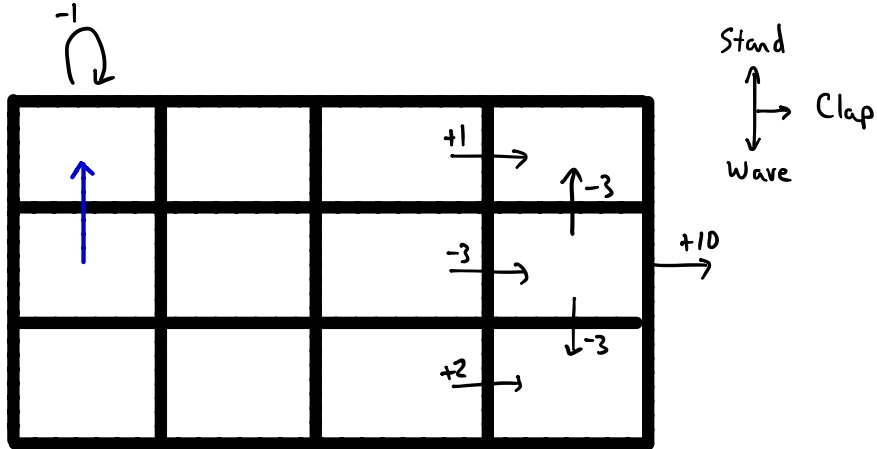| 0 | 0 | ~10.5 | ~9.5 |
| ~8.75 | ~8.75 | ~8.75 | ~10 |
| 0 | 0 | 0 | 0 |

| | | | |
| Δ0 | | | Δ0 |
| Δ0 | Δ0 | Δ2 | Δ10 |

— What about unvisited states?

— What if transition function were stochastic?

— Does the order of updates matter?

$\alpha = .5, \gamma = 1$
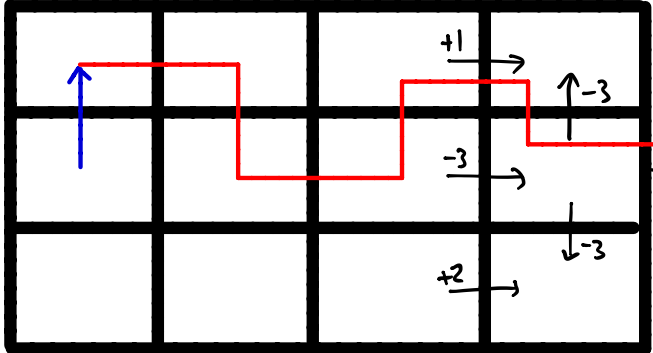
Stand
Clap
Wave

$S, -1$
$C, +1$
$W, 0$
$C, +10$
$C, 0$
$C, -3$
$C, +2$
$S, 0$

| 0 | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

← on policy

| 0 | 0 | ~10.5 | ~9.5 |
|---|---|---|---|
| ~8.75 | ~8.75 | ~8.75 | ~10 |
| 0 | 0 | 0 | 0 |

| | | | |
|---|---|---|---|
| Δ0 | | | Δ0 |
| Δ0 | Δ0 | Δ2 | Δ10 |

↑ then here etc.

← first updates

— What about unvisited states?

— What if transition function were stochastic?

— Does the order of updates matter?

  — prioritized sweeping

  — could actually do update in step (e)

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

Trajectory sampling: how differs from Dyna here?

$\alpha = .5, \gamma = 1$

Stand

Clap

Wave

Trajectory sampling: how differs from Dyna here?

How does Dyna-$Q^+$ differ from Dyna here?

— relationship to UCB?

— In what way does Dyna-$Q^+$ violate the authors' principles?

— Could you accomplish something similar w/out changing updates?

(see ex. 8.4)

# MCTS: Monte Carlo Tree Search — Planning at decision time

MCTS

Stand
Clap
Wave

-1

+1
-3
-3
+10   +11
-3
+2

MCTS

-1

+1

-3

-3

+10 +11

+4

+2

-3

Stand

Clap

Wave

MCTS

Stand
↑
← Clap
↓
Wave

# MCTS: Monte Carlo Tree Search – Planning at decision time



Stand
↑
├→ Clap
↓
Wave

- Interleaving planning and acting: model known
- Focusses search on current state
- Can combine w/ learning a model
- Can combine w/ a learned value function
- Random rollouts especially useful in game playing
- Can use more informed rollouts

# Approximation

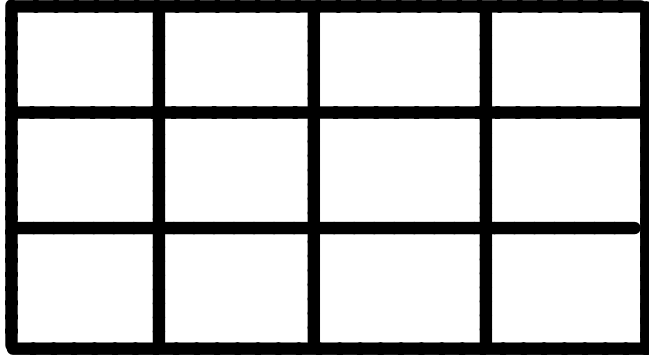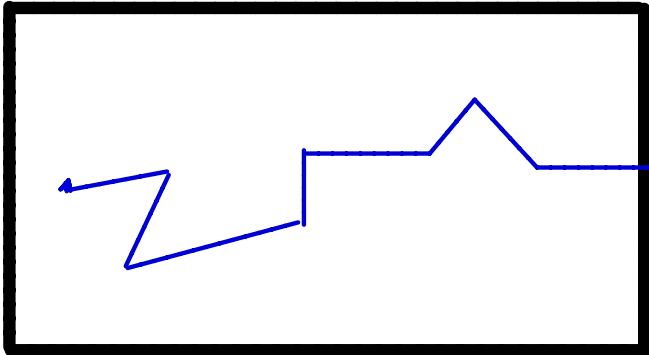|  |  |  |  |
|--|--|--|--|
|  |  |  |  |
|  |  |  |  |
|  |  |  |  |

Stand

Clap

Wave

# Approximation



Stand

Clap

Wave

Stand

Clap

Wave