Game tree search

$S_t$

$a_t$

$r_t$

$S_{t+1}$

$a_{t+1}$

Game tree search

Dynamic Programming (DP)

$S_t$

$a_t$

$r_t$

$S_{t+1}$

$a_{t+1}$

Game tree search

$S_t$

$a_t$

$r_t$

$S_{t+1}$

$a_{t+1}$

Dynamic Programming (DP)

Monte Carlo (MC)

Game tree search

Dynamic Programming (DP)

Monte Carlo (MC)

Temporal Difference (TD)

$S_t$

$a_t$

$r_t$

$S_{t+1}$

$a_{t+1}$

-1 ↻

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

Stand ↑
← Clap →
Wave ↓

| O | O | O | +1 → O |
| O | O | O | -3 → O |  +10 →
| O | O | O | +2 → O |

(with -3 ↑ between top right cells, -3 ↓ between middle right cells)

$$v(s) \leftarrow v(s) + \alpha [R + \gamma v(s') - v(s)]$$

"target"

$-1$



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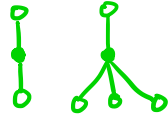
Stand

Clap

Wave

Grid (top-left) with cells:

| 0 | 0 | 0 | +1 → 0 |
| 0 | 0 | -3 → 0 | ↑ -3 → +10 |
| 0 | 0 | 0 +2 → | 0 |

↑ -3

↓ -3

Top-right grid:

| | 0 | | |
| 0 | 0 | | |
| | | | |

$v(s) \leftarrow v(s) + \alpha [\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$

-1

α=.5, γ=1

Stand
↑
→ Clap
↓
Wave

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

+10 →

| | 0 | ? | |
| 0 | 0 | | |
| | | | |

$$v(s) \leftarrow v(s) + \alpha [R + \gamma v(s') - v(s)]$$

"target"

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

-1

Stand
Clap
Wave

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

| | 0 | .5 | ? |
|---|---|---|---|
| 0 | 0 | | ? |
| | | | |

$$v(s) \leftarrow v(s) + \alpha [R + \gamma v(s') - v(s)]$$

"target"

-1

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

Stand
↑
→ Clap
↓
Wave

Grid (top-left) cells all contain 0, with transitions:
+1, -3, +10, -3, -3, +2

Grid (top-right, blue):

| | 0 | .5 | -25 |
|---|---|---|---|
| 0 | 0 | | 5 |
| | | | |

| | | | |
|---|---|---|---|
| 0 | 0 | ? | ? |
| | | | |

$v(s) \leftarrow v(s) + \alpha [\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$

α=.5,  γ=1

-1 ↺

Stand ↑ → Clap
Wave ↓

Grid (3×4):

| 0 | 0 | 0 +1→ | 0 |
|---|---|---|---|
| 0 | 0 | 0 -3→ | 0 ↑-3 +10→ |
| 0 | 0 | 0 +2→ | 0 ↓-3 |

Right tables:

|   | 0 | .5 | -.25 |
|---|---|----|------|
| 0 | 0 |    | 5 |
|   |   |    |   |

|   |   |   |     |
|---|---|---|-----|
| 0 | 0 | 1 | 7.5 |
|   |   |   |     |

$$v(s) \leftarrow v(s) + \alpha[\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$$

-1

α=.5, γ=1

Stand
↑
→ Clap
↓
Wave



Grid world with cells containing O (zeros), with transitions labeled: +1, -3, -3, +10, -3, +2

| | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | | 5 |
| | | | |

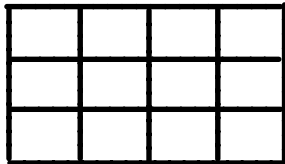| 0 | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

| | | ? | ? |
|---|---|---|---|
| 0 | .5 | .75 | 8.75 |
| | | | |

$$v(s) \leftarrow v(s) + \alpha [\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$$

-1

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

Stand
↑
→ Clap
↓
Wave

TD

| | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | | 5 |
| | | | |

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

| 0 | 0 | .625 | 3.625 |
|---|---|---|---|
| 0 | .5 | .75 | 8.75 |
| 0 | 0 | 0 | 0 |

Grid (top-left):

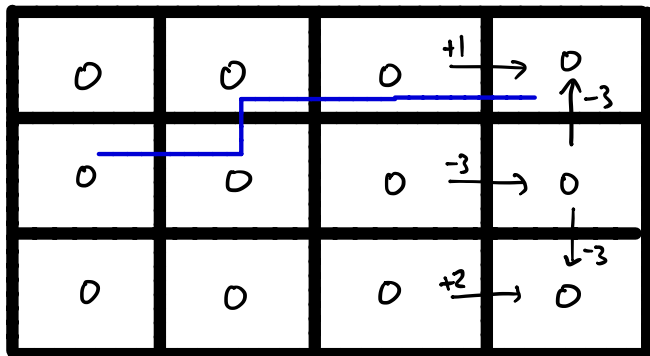| 0 | 0 | 0 +1→ | 0 |
|---|---|---|---|
| 0 | 0 | 0 -3→ | 0 -3↑ |
| 0 | 0 | 0 +2→ | 0 -3↓ |

+10

$$v(s) \leftarrow v(s) + \alpha \underbrace{[R + \gamma v(s') - v(s)]}_{\text{"target"}}$$

MC (first visit)

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

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

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

-1

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

Stand ↑ → Clap
Wave ↓

TD

+1 →
-3 ↑
-3 →
+10 →
-3 ↓
+2 →

| | O | O | O |
|---|---|---|---|
| O | O | O | O |
| O | O | O | O |

(grid values, all O)

TD tables:

| | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | | 5 |
| | | | |

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

| 0 | 0 | .625 | 3.625 |
|---|---|---|---|
| 0 | .5 | .75 | 8.75 |
| 0 | 0 | 0 | 0 |

$v(s) \leftarrow v(s) + \alpha [\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$

MC (first visit)

| | 5 | 5 | 4.5 |
|---|---|---|---|
| 5 | 5 | | 5 |
| | | | |

| | | | |
|---|---|---|---|
| ? | ? | ? | 7.5 |
| | | | |

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

-1 ↺

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

Legend: Stand ↑ / Clap → / Wave ↓

Grid rewards: +1, -3, +10, -3, -3, +2

**TD**

|     | 0   | .5  | -.25 |
| --- | --- | --- | ---- |
| 0   | 0   |     | 5    |
|     |     |     |      |

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

| 0   | 0   | .625 | 3.625 |
| --- | --- | ---- | ----- |
| 0   | .5  | .75  | 8.75  |
| 0   | 0   | 0    | 0     |

$$v(s) \leftarrow v(s) + \alpha \underbrace{[R + \gamma v(s') - v(s)]}_{\text{"target"}}$$

**MC (first visit)**

|     | 5   | 5   | 4.5 |
| --- | --- | --- | --- |
| 5   | 5   |     | 5   |
|     |     |     |     |

| 0   | 5   | 5   | 4.5 |
| --- | --- | --- | --- |
| 6   | 6   | 3.5 | 7.5 |
| 0   | 0   | 0   | 0   |

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

-1 ↺

α=.5, γ=1

Stand ↑ / Clap → / Wave ↓



| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 (+1 →) |
| 0 | 0 | 0 (−3 →) | 0 (→) |
| 0 | 0 (+2 →) | 0 | 0 |

with rewards: +1, −3, +10, −3, −3, +2

**TD**

| | | | |
|---|---|---|---|
| | 0 | .5 | -.25 |
| 0 | 0 | | 5 |

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

| | | | |
|---|---|---|---|
| 0 | 0 | .625 | 3.625 |
| 0 | .5 | .75 | 8.75 |
| 0 | 0 | 0 | 0 |

$$v(s) \leftarrow v(s) + \alpha[\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$$

TD(0):
- one-step
- tabular
- model free

**MC (first visit)**

| | | | |
|---|---|---|---|
| | 5 | 5 | 4.5 |
| 5 | 5 | | 5 |
| | | | |

| | | | |
|---|---|---|---|
| 0 | 5 | 5 | 4.5 |
| 6 | 6 | 3.5 | 7.5 |
| 0 | 0 | 0 | 0 |

| | | | |
|---|---|---|---|
| | | 8 | 7.25 |
| 8.5 | 8.5 | 7.25 | 8.75 |
| | | | |

-1 ↻

α=.5, γ=1

Stand
↑
→ Clap
↓
Wave

**TD**

Gridworld:

| | | ③ | ② |
|---|---|---|---|
| O | O | O +1→ | O |
| ① O | O | ④ O -3→ | O -3↓ +10→ |
| O | O | O +2→ | O |

(arrows: -3↑, -3↓)

TD table 1:
| | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | | 5 |
| | | | |

TD table 2 (red):
| 0 | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | 1 | 7.5 |
| 0 | 0 | 0 | 0 |

TD table 3 (green):
| 0 | 0 | .625 | 3.625 |
|---|---|---|---|
| 0 | .5 | .75 | 8.75 |
| 0 | 0 | 0 | 0 |

$$v(s) \leftarrow v(s) + \alpha[\underbrace{R + \gamma v(s') - v(s)}_{\text{"target"}}]$$

**TD(0):**
- one-step
- tabular
- model free

**MC (first visit)**

MC table 1:
| | 5 | 5 | 4.5 |
|---|---|---|---|
| 5 | 5 | | 5 |
| | | | |

MC table 2 (red):
| 0 | 5 | 5 | 4.5 |
|---|---|---|---|
| 6 | 6 | 3.5 | 7.5 |
| 0 | 0 | 0 | 0 |

MC table 3 (green):
| | | 8 | 7.25 |
|---|---|---|---|
| 8.5 | 8.5 | 7.25 | 8.75 |
| | | | |

Batch MC
?

Batch TD
?

Convergence of Batch methods
① ② ③ ④

| | Batch MC | Batch TD |
|---|---|---|
| ① | | |
| ② | | ? |
| ③ | | ? |
| ④ | ? | ? |

$-1$ ↺

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

Stand
→ Clap
Wave

**TD**

③ ②   +1 →   $-3$
①   $-3$ →   +10 →
$-3 ↓$
+2 →
④

$$v(s) \leftarrow v(s) + \alpha [R + \gamma v(s') - v(s)]$$

$\underbrace{\phantom{R + \gamma v(s') - v(s)}}_{\text{"target"}}$

TD(0):
- one-step
- tabular
- model free

TD tables:

|   | 0 | .5 | -.25 |
|---|---|----|------|
| 0 | 0 |    | 5    |
|   |   |    |      |

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

| 0 | 0 | .625 | 3.625 |
|---|---|------|-------|
| 0 | .5 | .75 | 8.75 |
| 0 | 0 | 0   | 0     |

**MC (first visit)**

|   | 5 | 5 | 4.5 |
|---|---|---|-----|
| 5 | 5 |   | 5   |
|   |   |   |     |

| 0 | 5 | 5   | 4.5 |
|---|---|-----|-----|
| 6 | 6 | 3.5 | 7.5 |
| 0 | 0 | 0   | 0   |

|     |     | 8    | 7.25 |
|-----|-----|------|------|
| 8.5 | 8.5 | 7.25 | 8.75 |
|     |     |      |      |

**Convergence of Batch methods**

Batch MC          Batch TD
① avg(7,11,10)=9.33
②                    9.5
③                    10.5
④    ?               ?

-1 ↺

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

Stand ↑ → Clap
Wave ↓

# TD

Grid world (large top-left grid):



TD(o):
- one-step
- tabular
- model free

$$v(s) \leftarrow v(s) + \alpha[R + \gamma v(s') - v(s)]$$
$\underbrace{\phantom{R + \gamma v(s')}}_{\text{"target"}}$

TD tables:

| | 0 | .5 | -.25 |
|---|---|---|---|
| 0 | 0 | | 5 |
| | | | |

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

| 0 | 0 | .625 | 3.625 |
|---|---|---|---|
| 0 | .5 | .75 | 8.75 |
| 0 | 0 | 0 | 0 |

# MC (first visit)

| | 5 | 5 | 4.5 |
|---|---|---|---|
| 5 | 5 | | 5 |
| | | | |

| 0 | 5 | 5 | 4.5 |
|---|---|---|---|
| 6 | 6 | 3.5 | 7.5 |
| 0 | 0 | 0 | 0 |

| | | 8 | 7.25 |
|---|---|---|---|
| 8.5 | 8.5 | 7.25 | 8.75 |
| | | | |

Convergence of Batch methods

Batch MC          Batch TD

① avg(7,11,10) = 9.33

②                              9.5

③                              10.5

④ avg(11,7) = 9     avg(10.5,7) = 8.75

$\alpha = .1 \quad \gamma = 1$

-1

Stand / Clap / Wave (TD: s a r s)

SARSA
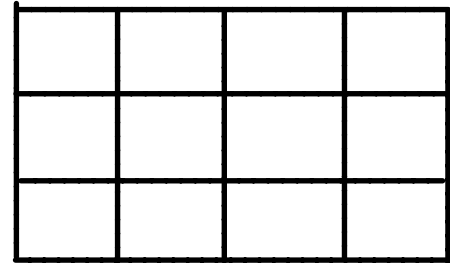
Policy: $\epsilon$-greedy, $\epsilon = .75$; Ties: →, ↓
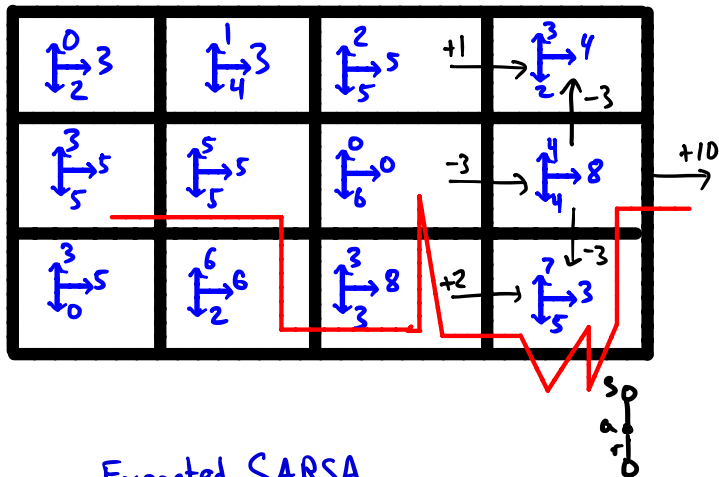
Expected SARSA

Q-learning

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ Target - Q(S,A) \right]$$

$\alpha = .1 \quad \gamma = 1$

-1

Stand
↑
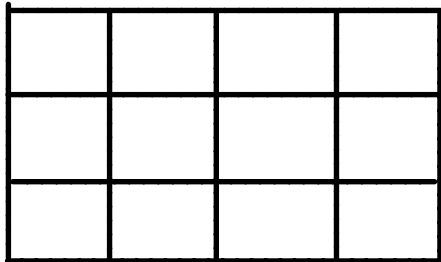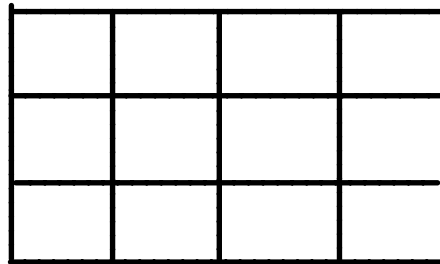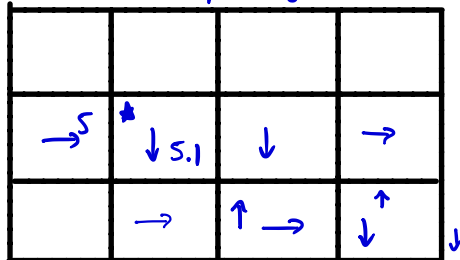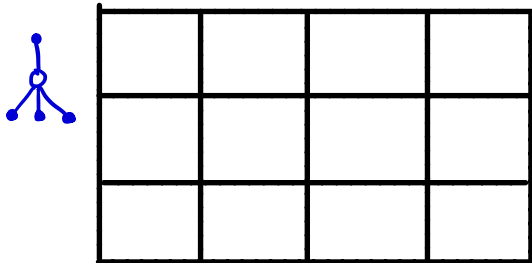↓ → Clap
Wave

**SARSA**

Policy: ε-greedy, ε = .75 ; Ties: →, ↓

→ 5    ↓?

**Expected SARSA**

**Q-learning**

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ \text{Target} - Q(S,A) \right]$$

Grid (top-left):

Row 1: cell 0: 3, 2 ; cell 1: 3, 4 ; cell 2: 5, 5 ; +1 → cell 3: 4, 2 ; -3 ↑
Row 2: cell 3: 5, 5 ; cell 5: 5, 5 ; cell 0: 0, 6 ; -3 → cell 4: 8, 4 ; +10 →
Row 3: cell 3: 5, 0 ; cell 6: 6, 2 ; +2 → cell 3: 8, 3 ; cell 7: 3, 5 ; -3

sp
a
r
b

$\alpha = .1$ $\gamma = 1$

-1

Stand → Clap
Wave

**SARSA**

☆ = policy change

→ 5 | ☆ ↓ 5.1 | ↓ | →
| → | ↑ → | ↑ ↓ | ↓

Policy: $\epsilon$-greedy, $\epsilon = .75$ ; Ties: →, ↓

**Expected SARSA**

**Q-learning**

+1, -3, -3, +10, +2, -3

0 3 / 2, 1 3 / 4, 2 5 / 5, 3 4 / 2
3 5 / 5, 5 5 / 5, 0 0 / 6, 4 8 / 4
3 5 / 0, 6 6 / 2, 3 8 / 3, 7 3 / 5

sp a r 0

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ \text{Target} - Q(S,A) \right]$$

-1

$\alpha = .1$  $\gamma = 1$

Stand
↑
↓ → Clap
Wave

Grid (top left):

| 0 →3 ↓2 | 1 →3 ↓4 | 2 →5 ↓5 | +1 → 3 →4 ↓2 ↑-3 |
| 3 →5 ↓5 | 5 →5 ↓5 | 0 →0 ↓6 | -3 → 4 →8 ↓4 | +10 → |
| 3 →5 ↓0 | 6 →6 ↓2 | +2 → 3 →8 ↓3 | 7 ↑-3 →3 ↓5 |

Sorb (so a r b)

## SARSA

★ = policy change

| | | | |
| →5 | ★ ↓5.1 | ↓6.2 | →8.2 |
| | ★ 5.7 → | ↑33 →7.9 | ↑7.1 ↓5.01  ↓ 4.9→5.01 |

Policy: ε-greedy, ε = .75 ; Ties: →, ↓

## Expected SARSA

(empty grid)

## Q-learning

| | | | |
| →5 | ★ ↓5.1 | | |
| | → | | |

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ \text{Target} - Q(S,A) \right]$$

-1

$\alpha = .1 \quad \gamma = 1$

Stand ↑ → Clap
Wave ↓

Main grid (blue):

| | | | |
|---|---|---|---|
| 0 →3 ↓2 | 1 →3 ↓4 | 2 →5 ↓5 | 3 →4 ↓2 |
| 3 →5 ↓5 | 5 →5 ↓5 | 0 →0 ↓6 | 4 →8 ↓4 |
| 3 →5 ↓0 | 6 →6 ↓2 | 3 →8 ↓3 | 7 →3 ↓5 |

+1 → ... -3 ... -3 → ... +10 → ... -3 ... +2 → ...

sarb

**SARSA**

★ = policy change

| | | | |
|---|---|---|---|
| | | | |
| →5 | ★ ↓5.1 | ↓6.2 | →8.2 |
| | ★ →5.7 | ↑33 →7.9 | ↑7.1 ↓5.01 |

↓4.9→5.01

Policy: ε-greedy, ε = .75 ; Ties: →, ↓

**Expected SARSA**

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

**Q-learning**

| | | | |
|---|---|---|---|
| | | | |
| →5 | ★ ↓5.1 | ↓6.2 | →8.2 |
| | →6.2 | ↑3.3 → | ↓ ↑7.1 |

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ Target - Q(S,A) \right]$$
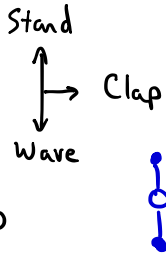
$\alpha = .1 \quad \gamma = 1$

Stand
↑
↓ → Clap
Wave

Grid (left, main):

| | | | |
|---|---|---|---|
| 0 →3 ↓2 | 1 →3 ↓4 | 2 →5 ↓5  +1→ | 3 →4 ↑2  -3 |
| 3 →5 ↓5 | 5 →5 ↓5 | 0 →0 ↓6  -3→ | 4 →8 ↓4  +10→ |
| 3 →5 ↓0 | 6 →6 ↓2 | 3 →8 ↓3  +2→ | 7 →3 ↑5  -3 |

**SARSA**

☆ = policy change

| | | | |
|---|---|---|---|
| | | | |
| →5 | ☆ ↓5.1 | ↓6.2 | →8.2 |
| | ☆ 5.7 → | ↑³³ →7.9 | ↑7.1 ↓5.01  ↓4.9→5.01 |

Policy: ε-greedy, ε = .75 ; Ties: →, ↓

**Expected SARSA**

| | | | |
|---|---|---|---|
| | | | |
| →5 | ↓ | | |
| | → | | |

**Q-learning**

| | | | |
|---|---|---|---|
| | | | |
| →5 | ☆ ↓5.1 | ↓6.2 | 8.2 → |
| | 6.2 → | ↑³·³ 8.1 → | ↑7.1 ↓5.19  ↓5.1→5.19 |

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ \text{Target} - Q(S,A) \right]$$

-1

$\alpha = .1 \quad \gamma = 1$

Stand
↑
|
→ Clap
|
↓
Wave

Main grid (top-left):
- Cell 0: →3, ↓2
- Cell 1: →3, ↓4
- Cell 2: →5, ↓5
- Cell 3: →4, ↑2, ↓3
- +1 →
- Cell 3: →5, ↓5
- Cell 5: →5, ↓5
- Cell 0: →0, ↓6
- Cell 4: →8, ↓4
- -3 →
- +10 →
- Cell 3: →5, ↓0
- Cell 6: →6, ↓2
- Cell 3: →8, ↓3
- Cell 7: →3, ↓5
- -3
- +2 →

S, a, r, b

## SARSA

☆ = policy change

| | | | |
|---|---|---|---|
| →5 | ☆ ↓5.1 | ↓6.2 | →8.2 |
| | ☆ 5.7 → | ↑33 →7.9 | ↑7.1 ↓5.01 |

↓4.9 → 5.01

Policy: ε-greedy, ε = .75 ; Ties: →, ↓

## Expected SARSA

| | | | |
|---|---|---|---|
| | | | |
| →5 | ↓5 | ↓5.95 | 8.2 → |
| | ☆ 5.95 → | ↑3 7.95 → | ↑6.9 ↓4.90375 |

↓4.95 → 4.90375

## Q-learning

| | | | |
|---|---|---|---|
| | | | |
| →5 | ☆ ↓5.1 | ↓6.2 | 8.2 → |
| | 6.2 → | ↑3.3 8.1 → | ↑7.1 ↓5.19 |

↓5.1 → 5.19

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ \text{Target} - Q(S,A) \right]$$

-1

$\alpha = .1 \quad \gamma = 1$

Stand
↑
→ Clap
↓
Wave

Main grid (4×3):

| ↑↓→ 0,3,2 | ↑↓→ 1,3,4 | ↑↓→ 2,5,5 | +1 → ↑↓→ 3,4,2 |
| ↑↓→ 3,5,5 | ↑↓→ 5,5,5 | ↑↓→ 0,0,6 | -3 → ↑↓→ 4,8,4  +10 → |
| ↑↓→ 3,5,0 | ↑↓→ 6,6,2 | ↑↓→ 3,8,3 | +2 → ↑↓→ 7,3,5  -3 |

↑ -3 (near top right)
-3 (middle)
reward
Sarbo

SARSA ← on policy

★ = policy change

|  |  |  |  |
|---|---|---|---|
|  |  |  |  |
| → 5 | ★ ↓ 5.1 | ↓ 6.2 | → 8.2 |
|  | ★ ↓→ 5.7 | ↑³³→ 7.9 | ↑ 7.1 ↓ 5.01  ↓ 4.9→5.01 |

Policy: ε-greedy, ε = .75 ; Ties: →, ↓

Expected SARSA ← can be off policy

|  |  |  |  |
|---|---|---|---|
|  |  |  |  |
| → 5 | ↓ 5 | ↓ 5.95 | → 8.2 |
|  | ★ → 5.95 | ↑³ → 7.95 | ↑ 6.9 ↓ 4.90375  ↓ 4.95→4.90375 |

Q-learning ← off policy

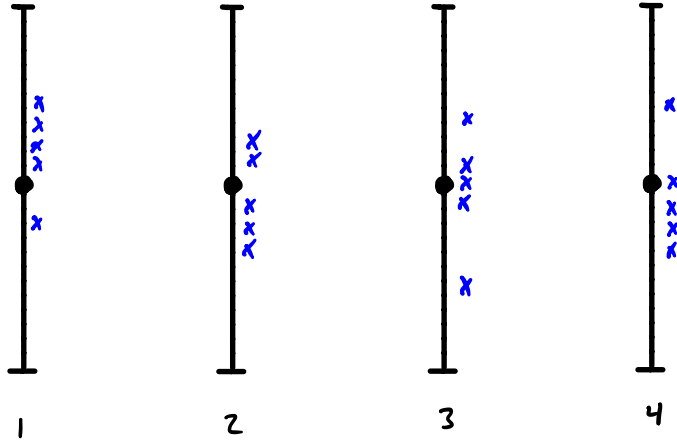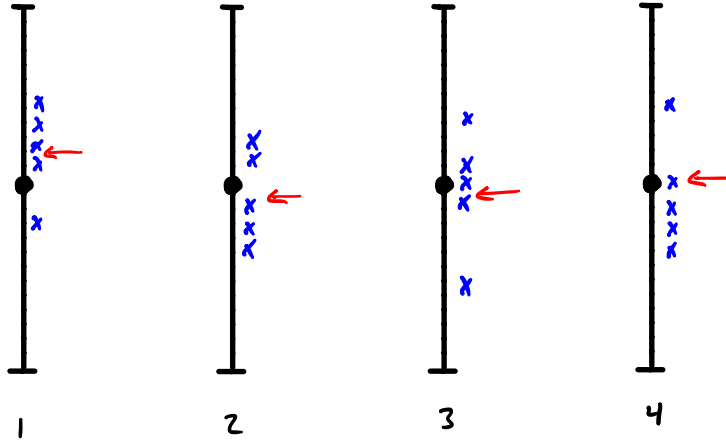|  |  |  |  |
|---|---|---|---|
|  |  |  |  |
| → 5 | ★ ↓ 5.1 | ↓ 6.2 | → 8.2 |
|  | → 6.2 | ↑³·³ → 8.1 | ↑ 7.1 ↓ 5.19  ↓ 5.1→5.19 |

- How do learned policies differ?
- Conditions for convergence?
- Why Expected SARSA not as known?

$$Q(S,A) \leftarrow Q(S,A) + \alpha \left[ \text{Target} - Q(S,A) \right]$$

Double Q learning — addresses maximization bias, regression to the mean
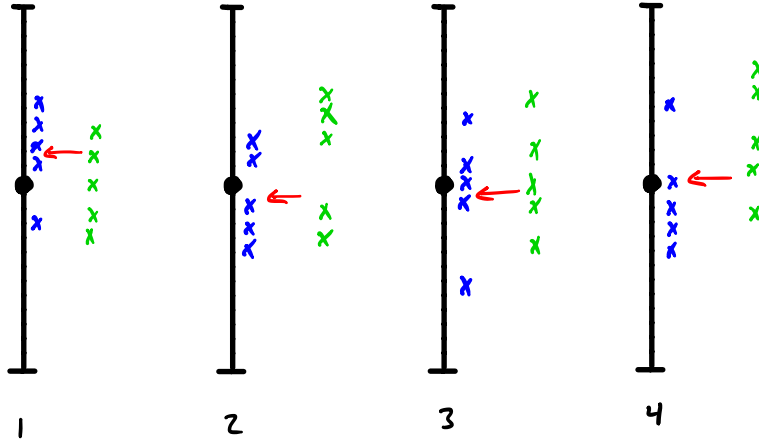(illustrated in bandit setting)

Double Q learning — addresses maximization bias, regression to the mean
(illustrated in bandit setting)



- true means
← sample means

1    2    3    4

Double Q learning — addresses maximization bias, regression to the mean
(illustrated in bandit setting)



• true means
← sample means
x new, independent samples

1    2    3    4

Ch 6    summary

Prediction:   TD(0)  = one-step, tabular, model-free TD  ⎤
                                                          ⎥  The core:
                                                          ⎥    bootstrapping
Control:  SARSA        ⎫ expected SARSA                   ⎥
          Q-learning   ⎭                                  ⎦


Also: Double Q
      After states