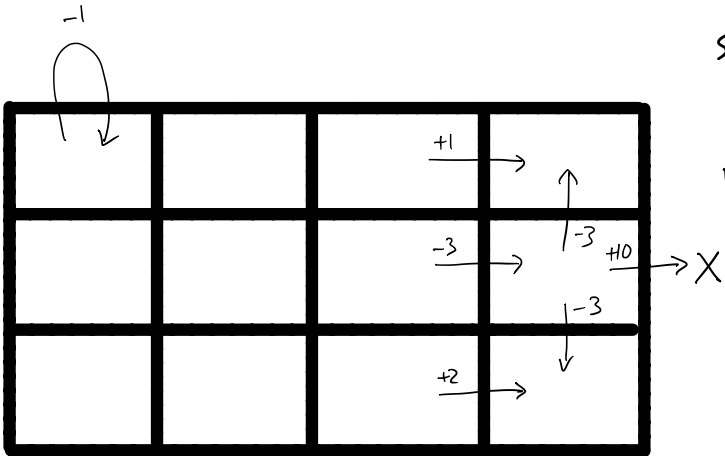Stand

Clap

Wave

(RL: no model)

Stand
↑
↔ Clap
↓
Wave

(DP: Known model)

-1

+1 →

-3 →  ↑  -3  +0 → X

-3 ↓

+2 →

# Policy Evaluation



Stand

Clap

Wave

# Policy Evaluation
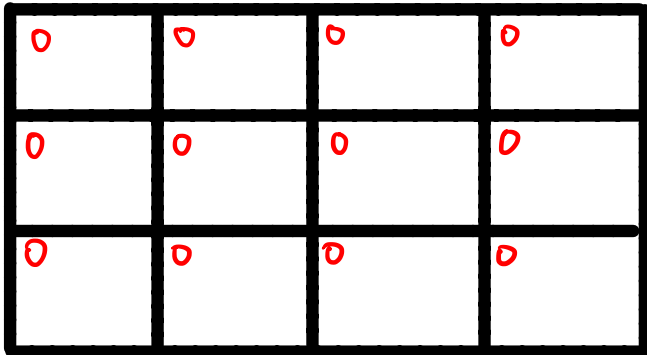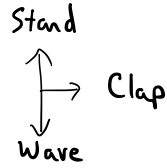


Stand ↑
→ Clap
Wave ↓

$\forall s \quad V_0(s) = 0$

$V_{k+1}(s) =$ (Eq. 4.5)

required

# Policy Evaluation

-1

Stand
Clap
Wave

+1

-3

-3 +0 → X

-3

+2

$\forall s \quad v_0(s) = 0$

$$v_{k+1}(s) = \sum_a \pi(a|s) \times \,?$$

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

# Policy Evaluation
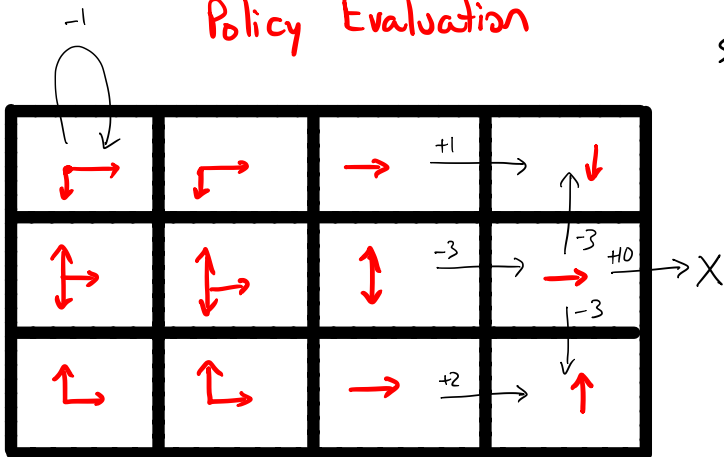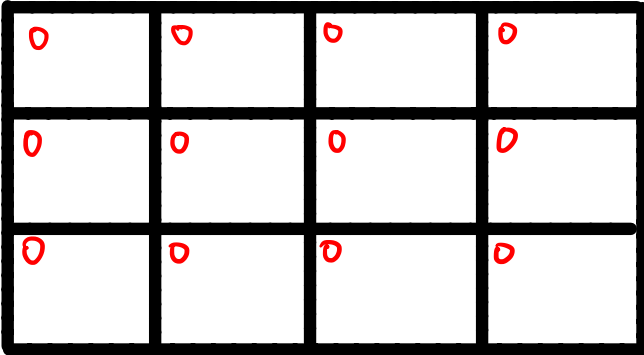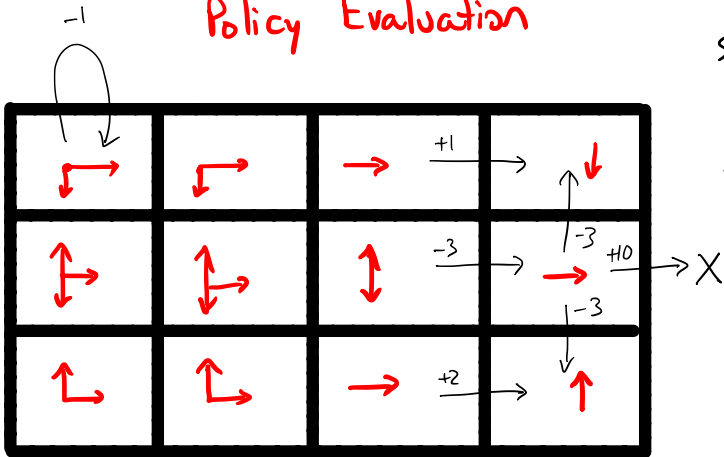


$$\forall s \quad v_0(s) = 0$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \times ?$$

# Policy Evaluation



Stand ← Clap (with Wave below, arrows indicating directions)
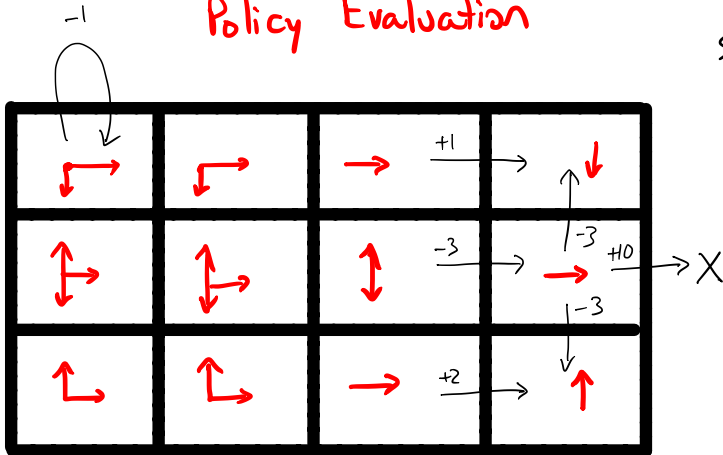
$$\forall s \quad v_0(s) = 0$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\left[r + \gamma v_k(s')\right]$$

# Policy Evaluation

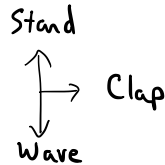

Stand
↑
→ Clap
↓
Wave

$$\forall s \quad v_0(s) = 0$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma V_k(s') \right]$$

| 0,0 | 0,0 | 0,1 | 0,0 |
| 0,0 | 0,0 | 0,0 | 0,10 |
| 0,0 | 0,0 | 0,2 | 0,0 |

# Policy Evaluation



Stand

Clap

Wave

$$\forall s \quad v_0(s) = 0$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\left[r + \gamma v_k(s')\right]$$

| | | | |
|---|---|---|---|
| 0,0,0 | 0,0,1,.5 | 0,1,1,11 | 0,0,10 |
| 0,0,0 | 0,0,0 | 0,0,1.5 ? | 0,10,10, |
| 0,0,0 | 0,0,1,1 | 0,2,2,12 | 0,0,10 |

# Policy Evaluation



Stand
Clap
Wave

$\forall s \quad v_0(s) = 0$

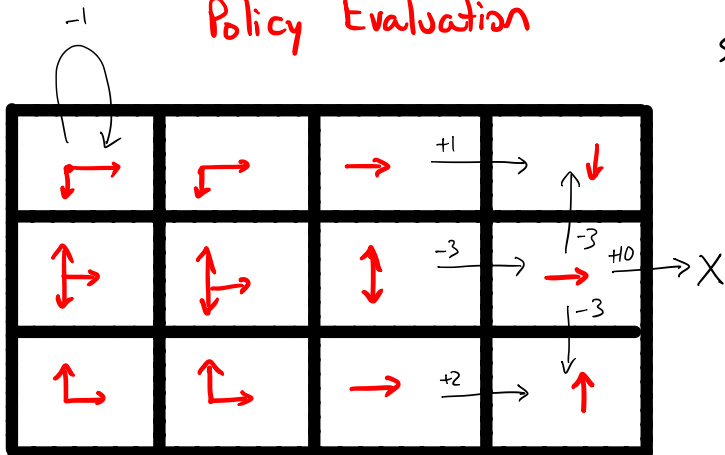$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')]$$

| 0,0,0 | 0,0,1,.5 | 0,1,1,11 | 0,0,10 |
| | ? | 11 | 10 |
| 0,0,0 | 0,0,0 | 0,0,15 | 0,10,10, |
| | | 11.5 | 10 |
| 0,0,0 | 0,0,1,1 | 0,2,2,12 | 0,0,10 |
| | | 12 | 10 |

# Policy Evaluation



Stand
↑
→ Clap
↓
Wave

$$\forall s \quad v_0(s) = 0$$

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma V_k(s') \right]$$

| | | | |
|---|---|---|---|
| 0,0,0<br>‖ 3/8 | 0,0,1,.5<br>‖ 1/4 | 0,1,1,1<br>‖ | 0,0,10<br>10 |
| 0,0,0<br>‖ 1/2 | 0,0,0<br>‖ 1/2 | 0,0,1.5<br>11.5 | 0,10,10,<br>10 |
| 0,0,0<br>‖ 3/8 | 0,0,1,1<br>‖ 3/4 | 0,2,2,12<br>12 | 0,0,10<br>10 |

Policy Improvement

Stand / Clap / Wave

Policy Improvement

Stand ↑
← → Clap
Wave ↓

X

Policy Improvement

Stand
↑
← → Clap
↓
Wave

Policy Improvement

Policy Improvement

Stand ↑
Clap →
Wave ↓

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

Left column value tables:

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

| 0 | 0 | 1 | 0 |
|---|---|---|----|
| 0 | 0 | 0 | 10 |
| 0 | 0 | 2 | 0 |

| 0 | 1 | 1 | 10 |
|---|---|---|----|
| 0 | 0 | 7 | 10 |
| 0 | 2 | 2 | 10 |

| 1 | 1 | 11 | 10 |
|---|---|----|----|
| 0 | 7 | 7  | 10 |
| 2 | 2 | 12 | 10 |

Right column (Policy Improvement) tables — arrows and a "?"

# Policy Improvement

Stand ↑
Clap →
Wave ↓

−1

+1

Stand → Clap, Wave

−3

−3  +0 → X

−3

+2

**Table 1 (values):**

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

**Table 1 (policy):**

| ↓→ | ↱ | → | ↓ |
| ↕→ | ↕→ | ↓ | → |
| ↱ | ↱ | → | ↑ |

**Table 2 (values):**

| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 10 |
| 0 | 0 | 2 | 0 |

**Table 2 (policy):**

| | → | | |
| | | → | |
| | → | | |

**Table 3 (values):**

| 0 | 1 | 1 | 10 |
| 0 | 0 | 7 | 10 |
| 0 | 2 | 2 | 10 |

**Table 3 (policy):**

| → | → | → | ↓ |
| ↕→ | → | → | ↑ |
| → | → | → | ↑ |

**Table 4 (values):**

| 1 | 1 | 11 | 10 |
| 0 | 7 | 7 | 10 |
| 2 | 2 | 12 | 10 |

**Table 4 (policy):**

| → | → | → | ↓ |
| → | → | ↓ | → |
| → | → | → | ↑ |

**Table 5 (values):**

| | | | |
| | | ? | |
| | ? | | |

**Table 5 (policy):** (empty)

**Table 6 (values):** (empty)

**Table 6 (policy):** (empty)

**Policy Improvement**

-1

Stand
Clap
Wave

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

X

Policy tables (values):

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

| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 10 |
| 0 | 0 | 2 | 0 |

| 0 | 1 | 1 | 10 |
| 0 | 0 | 7 | 10 |
| 0 | 2 | 2 | 10 |

| 1 | 1 | 11 | 10 |
| 0 | 7 | 7 | 10 |
| 2 | 2 | 12 | 10 |

| 1 | 11 | 11 | 10 |
| 7 | 7 | 12 | 10 |
| 2 | 12 | 12 | 10 |

(empty grid)

Improvement (policy arrows):

Grid 1, Grid 2 (→), Grid 3, Grid 4, with "?" in last populated grid, empty grid last.

Policy Improvement

-1

+1

Stand

Clap

Wave

-3

-3

+0

X

-3

+2

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

| ↓→ | ↓→ | → | ↓ |
| → | ↓→ | ↓ | → |
| ↳ | ↳ | → | ↑ |

| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 10 |
| 0 | 0 | 2 | 0 |

| | → | | |
| | | → | |
| | → | | |

| 0 | 1 | 1 | 10 |
| 0 | 0 | 7 | 10 |
| 0 | 2 | 2 | 10 |

| → | → | → | ↓ |
| ↳ | → | → | → |
| → | → | → | ↑ |

| 1 | 1 | 11 | 10 |
| 0 | 7 | 7 | 10 |
| 2 | 2 | 12 | 10 |

| → | → | → | ↓ |
| → | → | ↓ | → |
| → | → | → | ↑ |

| 1 | 11 | 11 | 10 |
| 7 | 7 | 12 | 10 |
| 2 | 12 | 12 | 10 |

| → | → | ↓ | ↓ |
| → | ↓→ | ↓ | → |
| → | → | ↑ | ↑ |
$= \pi^*$

| 11 | 11 | 12 | 10 |
| 7 | 12 | 12 | 10 |
| 12 | 12 | 12 | 10 |

| | | ↓ | |
| ↳ | ↓→ | ↓ | |
| | ↳ | ↳ | |
$= \pi^*$ (also)

Value Iteration

# Policy Improvement



-1

Stand
Clap
Wave

+1   -3   +0 → X   -3   +2

**Policy Iteration**

?

**Value Iteration**

| 0 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

| 0 | 0 | 1 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 10 |
| 0 | 0 | 2 | 0 |

| 0 | 1 | 1 | 10 |
|---|---|---|---|
| 0 | 0 | 7 | 10 |
| 0 | 2 | 2 | 10 |

| 1 | 1 | 11 | 10 |
|---|---|---|---|
| 0 | 7 | 7 | 10 |
| 2 | 2 | 12 | 10 |

| 1 | 11 | 11 | 10 |
|---|---|---|---|
| 7 | 7 | 12 | 10 |
| 2 | 12 | 12 | 10 |

| 11 | 11 | 12 | 10 |
|---|---|---|---|
| 7 | 12 | 12 | 10 |
| 12 | 12 | 12 | 10 |

$= \pi^*$

$= \pi^*$ (also)

# Policy Improvement

Stand ↑ → Clap
Wave ↓

-1

+1

-3   +0   → X
-3

+2

**Policy Iteration**

**Value Iteration**

-1

Stand
Clap
Wave

+1

-3
+10 → X
-3
-3

+2

# Policy Improvement

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

| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 10 |
| 0 | 0 | 2 | 0 |

| 0 | 1 | 1 | 10 |
| 0 | 0 | 7 | 10 |
| 0 | 2 | 2 | 10 |

| 1 | 1 | 11 | 10 |
| 0 | 7 | 7 | 10 |
| 2 | 2 | 12 | 10 |

| 1 | 11 | 11 | 10 |
| 7 | 7 | 12 | 10 |
| 2 | 12 | 12 | 10 |

| 11 | 11 | 12 | 10 |
| 7 | 12 | 12 | 10 |
| 12 | 12 | 12 | 10 |

$= \pi^*$

$= \pi^*$ (also)

# Policy Iteration

| 11 3/8 | 11 1/4 | 11 | 10 |
| 11 1/2 | 11 1/2 | 11 1/2 | 10 |
| 11 3/8 | 11 3/4 | 12 | 10 |

| 12 | 12 | 12 | 10 |
| 12 | 12 | 12 | 10 |
| 12 | 12 | 12 | 10 |

# Value Iteration