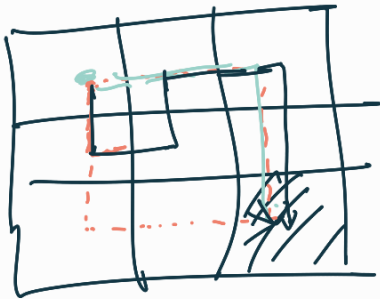# Dynamic Programming

→ policy evaluation
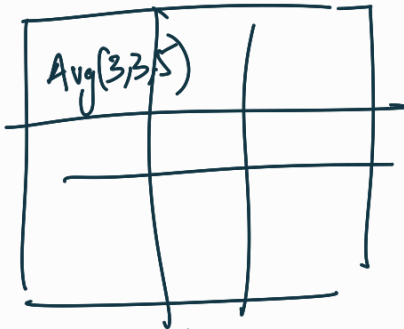→ policy improvement — Value iteration

$$\pi(a|s) \rightarrow policy$$

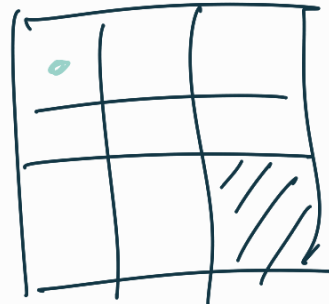R1  $\pi_1$ =



3, 3, 5



Avg(3,3,5)

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

| 4 | 5 | 3 |
|---|---|---|
| 0 | 2 | 2 |
| 1 | 1 | 0 |

9

R2  $\pi_2$



3,5,7



Avg(3,5,7)

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

| 5 | 10 | 12 |
|---|----|----|
| 14 | 13 | 3 |
| 2 | 5 | 0 |

R3  $\pi_3$



5, 10, 15



Avg
(5,10,15)

$$V_{\pi_3}(s) =$$

| 10 | 12 | 13 |
|----|----|----|
| 15 | 14 | 4 |
| 3 | 6 | 0 |

$$V_{\pi_3}(s) > V_{\pi_2}(s) > V_{\pi_1}(s)$$

# Policy Evaluation

$$v_\pi(s) = \sum_{a \in A} \pi(a|s) \left[ \sum_{s' \in S} P(s'|s,a)\left(r(s,a,s') + \gamma v_\pi(s')\right)\right]$$

?
0

Bellmans eqn.

?
0

init,
$$v_\pi(s) = 0 \leftarrow$$

Iterate:

$$v_{New}(s) = \sum_{a \in A} \pi(a|s) \left[ \sum_{s' \in S} P(s'|s,a)\left(r(s,a,s') + \gamma v_{OLD}(s')\right)\right]$$

$$v_{OLD}(s) = v_{New}(s)$$

-1

$\{(0,0),(0,1)(1,0)\ldots\}$  until $\|v_{OLD} - v_{New}\| \leq tol$

0.0001

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

$$diff_1 = 5$$

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

$v_{OLD}$

$$diff_2 = 3$$

$v_{New}$

2

# Policy improvement

### value iteration

$$\pi(a|s) \rightarrow \text{Not known}$$

we need to come up
with optimal ~~o~~ policy
$$\pi^*(a|s)$$

$$v_*(s) = \max_{a \in A} \left( \sum_{s' \in S} P(s'|s,a) \underbrace{\left( \underline{r(s,a,s') + \gamma v_*(s')} \right)}_{q(s,a)} \right)$$

Bellmans optimality equation

*  $$V_{OLD}(s) = 0$$

$$V_{NEW}(s) = \max_{a \in A} \left( \sum_{s' \in S} P(s'|s,a) \left( r(s,a,s') + \gamma V_{OLD}(s) \right) \right)$$

repeat until $\| V_{OLD}^{(s)} - V_{NEW}^{(s)} \| \leq \Delta$

$$Ex, \Delta = 0.001$$

Finally          after convergence

$$v_{NEW}(s) \approx v^*(s)$$

optimal state
value function

$v^*(s) \xrightarrow{\text{gives}} \pi^*(a|s)$ optimal policy

| 6 | 5 | 7 |
|---|---|---|
| 1 | 3 | 2 |
| 4 | 8 | 10 |

one-step reward $= -1$

$$5 - 1 = 4$$

$$0 \leftarrow \quad \rightarrow 1$$

8