# Week 8: MDP and Value Iteration

COMP90054 – AI Planning for Autonomy

# Key concepts

- Markov Decision Processes (MDPs)

- Solving MDPs:
  - Value Iteration
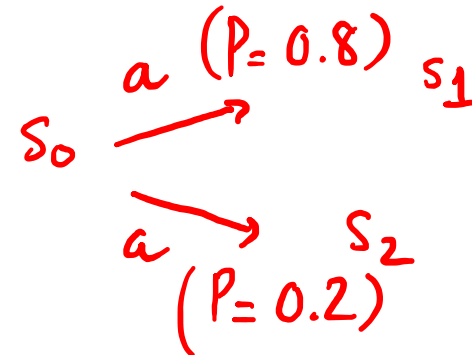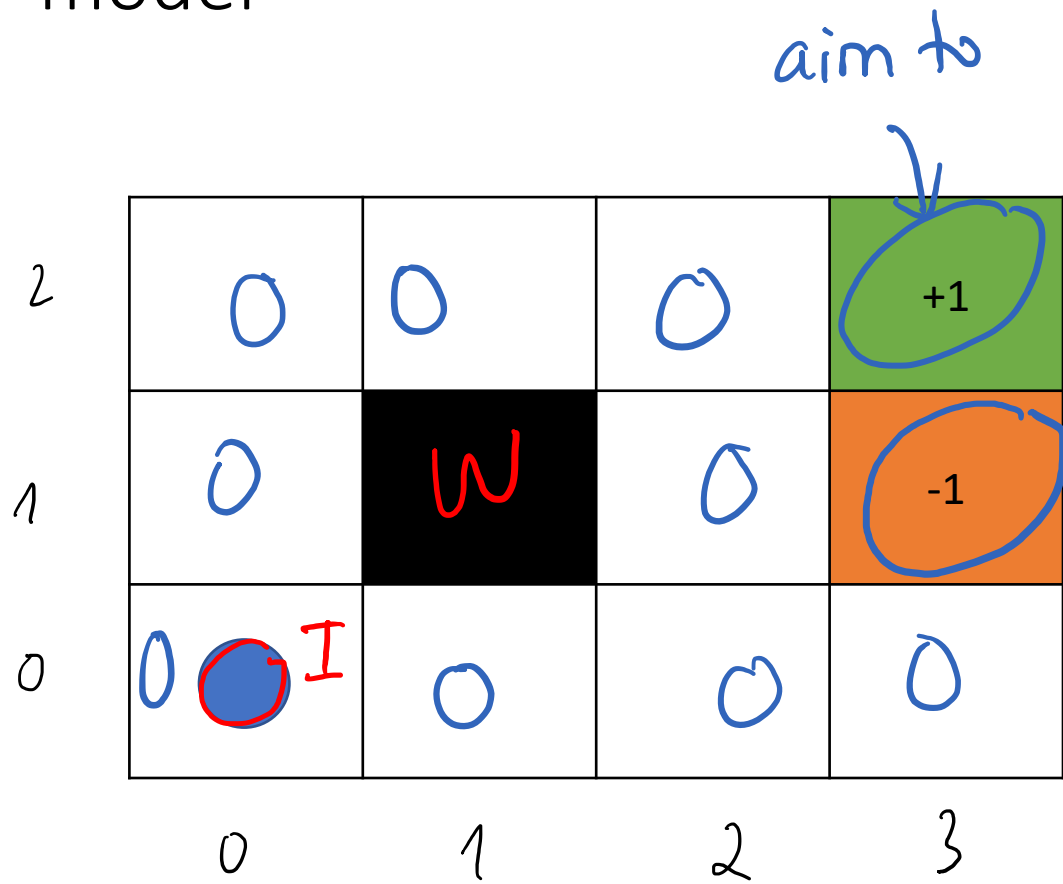
# Classical Planning vs. MDPs

*minimise the cost* (handwritten)

*maximise the reward* ✗ (handwritten)

| Classical Planning (6) | Markov Decision Processes (MDPs) |
|---|---|
| Set of states S | Set of states S |
| Initial state $s_0$ | Initial state $s_0$ |
| Action A(s) | Action A(s) |
| *Transition function s' = f(a, s)* | *Transition probabilities $P_a(s'\|s)$* Non-deterministic |
| *Goals $S_G \subseteq S$* | *Reward function r(s, a, s') (positive or negative)* |
| *Action costs c(a, s)* | |
| | Discount factor $0 \leq \gamma \leq 1$ (prefer shorter plans over longer plans) |

$$s_0 \xrightarrow{a} s_1$$

(single outcome)

$$s_0 \xrightarrow{a \ (P = 0.8)} s_1$$
$$\xrightarrow{a} s_2$$
$$(P = 0.2)$$

# Task 1: Model the Grid MDP example with a formal discounted-reward MDP model

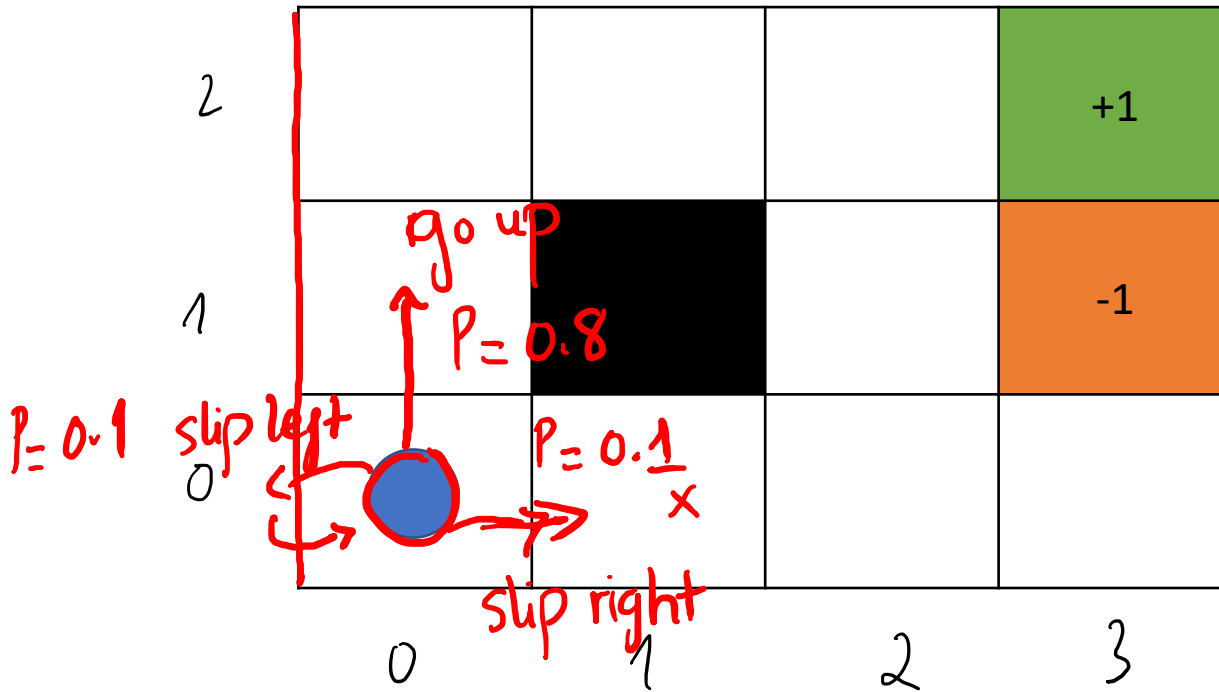$S, s_0, A(s), P_a(s'|s), r(s, a, s'), \gamma$



$$S = \{\langle x, y \rangle \mid x \in \{0, 1, 2, 3\},$$
$$y \in \{0, 1, 2\},$$
$$\langle 1, 1 \rangle \text{ is the wall}\}$$

$$s_0 = \langle 0, 0 \rangle$$
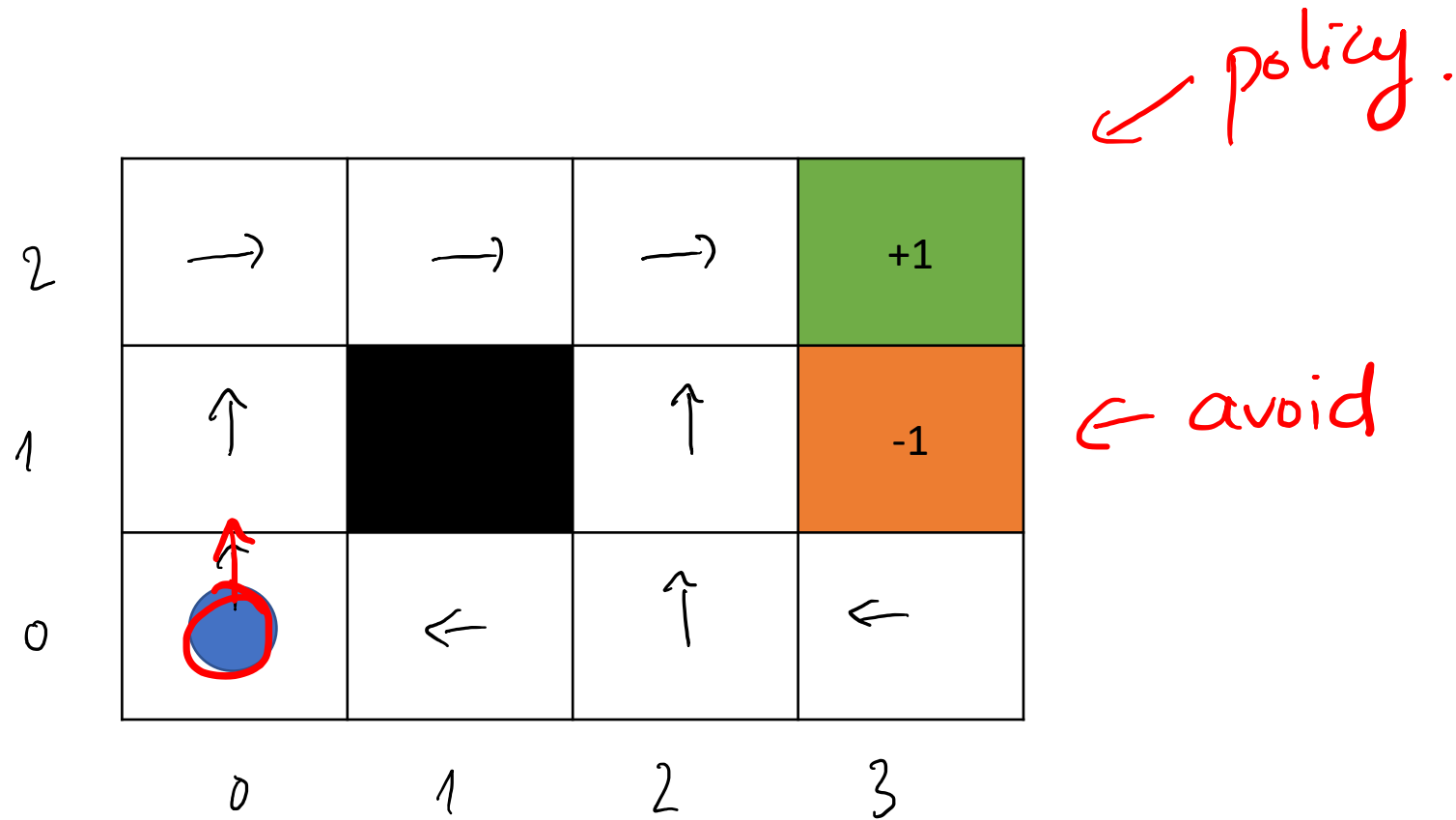
$$A = \{ \text{North, South, East, West} \}$$

# Task 1: Model the Grid MDP example with a formal discounted-reward MDP model

S, $s_0$, A(s), $P_a(s'|s)$, $r(s, a, s')$, $\gamma$



- Action: North:

+ P = 0.8 : go up.

+ P = 0.1: slip right

+ P = 0.1: slip left

- reward.?

# Task 1: Model the Grid MDP example with a formal discounted-reward MDP model

# Solving MDPs?

**Bellman equations**

_maximise the reward._

_get the expected reward of an action._

For discounted-reward MDPs the Bellman equation is defined recursively as:

expected
reward of action
a at state s

$$Q(s,a) = \sum_{s' \in S} P_a(s'|s) \left[ r(s,a,s') + \gamma V(s') \right]$$

Q-value

the probability
of action a

immediate
reward

discounted
future reward

$$V(s) = \max_{a \in A(s)} Q(s,a)$$

expected value of being in state s
and acting optimally

# Solving MDPs? Value Iteration

■ Set $V_0$ to arbitrary value function; e.g., $V_0(s) = 0$ for all $s$.

■ Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma \, V_i(s') \right]$$
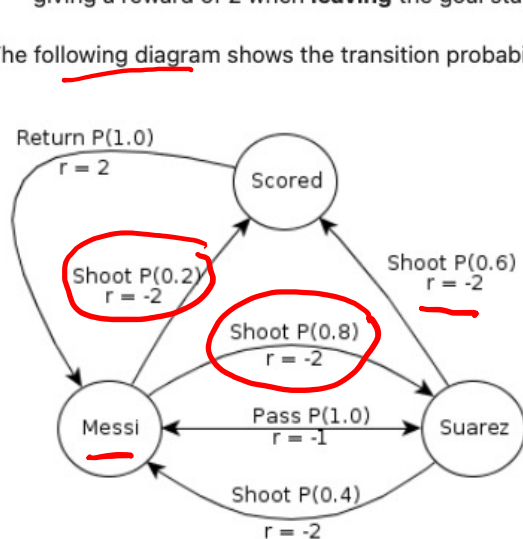
# Workshop Problem

Consider two football-playing robots: Messi and Suarez.

They play a simple two-player cooperate game of football, and you need to write a controller for them. Each player can pass the ball or can shoot at goal.

The football game can be modelled as a discounted-reward MDP with three states: *Messi*, *Suarez* (denoting who has the ball), and *Scored* (denoting that a goal has been scored); and the following action descriptions:

- If Messi shoots, he has 0.2 chance of scoring a goal and a 0.8 chance of the ball going to Suarez. Shooting towards the goal incurs a cost of 2 (or a reward of −2).

- If Suarez shoots, he has 0.6 chance of scoring a goal and a 0.4 chance of the ball going to Messi. Shooting towards the goal incurs a cost of 2 (or a reward of −2).

- If either player passes, the ball will reach its intended target with a probability of 1.0. Passing the ball incurs a cost 1 (or a reward of −1).

- If a goal is scored, the only action is to return the ball to Messi, which has a probability of 1.0 and has a reward of 2. Thus the reward for scoring is modelled by giving a reward of 2 when **leaving** the goal state.

The following diagram shows the transition probabilities and rewards:



shoot

Messi → Scored  $P = 0.2$

shoot → Suarez  $P = 0.8$

$P_{shoot}(Suarez \mid Messi) = 0.8$

to          from

$S = \{Messi, Suarez, Scored\}$

$r(Messi, shoot, Suarez) = -2$

Find the best action for all states?

# Workshop Problem

| | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| V(Messi) | 0 | | | |
| V(Suarez) | 0 | | | |
| V(Scored) | 0 | | | |

Iteration 0: Set $V_0(s) = 0$ for all s

The following diagram shows the transition probabilities and rewards:



Return P(1.0)
r = 2
Scored
Shoot P(0.2)
r = -2
Shoot P(0.6)
r = -2
Shoot P(0.8)
r = -2
Messi
Pass P(1.0)
r = -1
Suarez
Shoot P(0.4)
r = -2

$\gamma = 1$

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 |  |  |
| **V(Suarez)** | 0 |  |  |  |
| **V(Scored)** | 0 |  |  |  |

■ Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma \, V_i(s') \right]$$

The following diagram shows the transition probabilities and rewards:



s = Messi
s' = Suarez/Scored
a = shoot/pass
$\gamma = 1$

Iteration 1: $V_1(Messi)$

- shoot

$$P_{shoot}(Suarez|Messi)\left[r(Messi, shoot, Suarez) + \gamma \times V(Suarez)\right]$$
$$+ P_{shoot}(Scored|Messi)\left[r(Messi, shoot, Scored) + \gamma \times V(Scored)\right]$$
$$= 0.8\left[-2 + 1 \times 0\right] + 0.2\left[-2 + 1 \times 0\right] = -2$$

max → -1

- pass

$$P_{pass}(Suarez|Messi)\left[r(Messi, pass, Suarez) + \gamma \times V(Suarez)\right]$$
$$= 1 \times \left[-1 + 1 \times 0\right] = -1$$

Thao Le

11

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 |  |  |
| **V(Suarez)** | 0 | -1 |  |  |
| **V(Scored)** | 0 |  |  |  |

■ Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma V_i(s') \right]$$

Iteration 1: $V_1(Suarez)$

- shoot


- pass

The following diagram shows the transition probabilities and rewards:



s = Suarez
s' = Messi/Scored
a = shoot/pass
$\gamma = 1$

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 |  |  |
| **V(Suarez)** | 0 | -1 |  |  |
| **V(Scored)** | 0 | 2 |  |  |

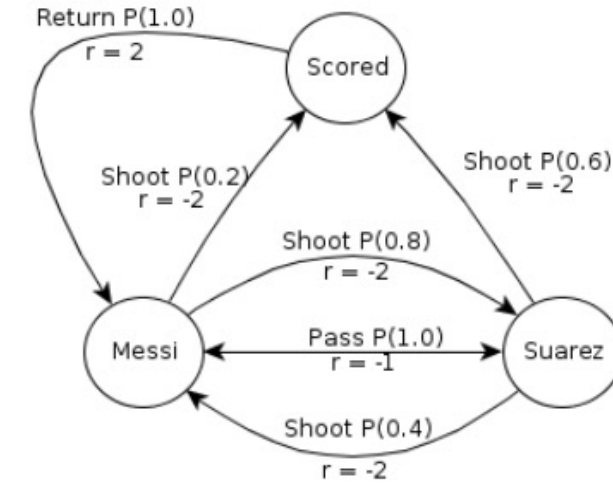

- Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma\, V_i(s') \right]$$

Iteration 1: $V_1(Scored)$

- return

s = Scored
s' = Messi
a = return
$\gamma = 1$

# Workshop Problem

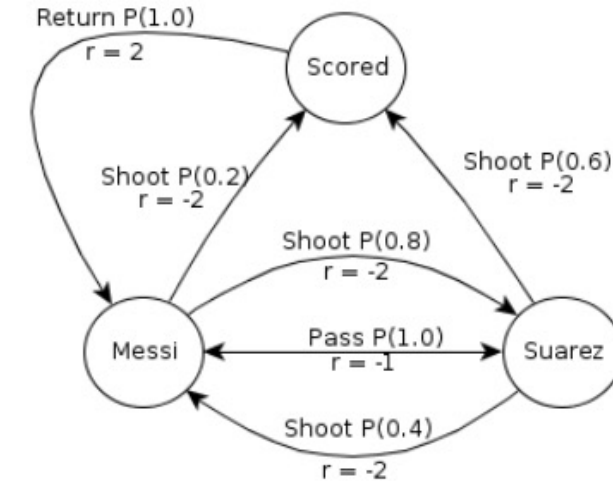|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 | -2 | |
| **V(Suarez)** | 0 | -1 | | |
| **V(Scored)** | 0 | 2 | | |



■ Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma \, V_i(s') \right]$$

Iteration 2: $V_2(Messi)$

- shoot

- pass

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 | -2 | |
| **V(Suarez)** | 0 | -1 | -1.2 | |
| **V(Scored)** | 0 | 2 | | |



Return P(1.0)
r = 2
Scored

Shoot P(0.2)
r = -2

Shoot P(0.6)
r = -2

Shoot P(0.8)
r = -2

Messi

Pass P(1.0)
r = -1

Suarez

Shoot P(0.4)
r = -2
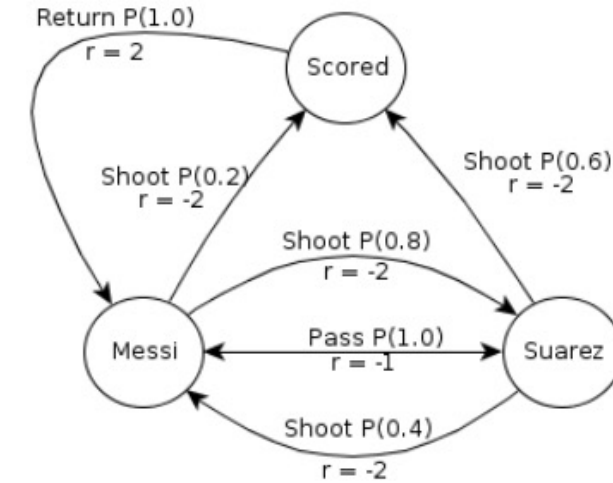
- Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma V_i(s') \right]$$

Iteration 2: $V_2(Suarez)$

- shoot

- pass

# Workshop Problem

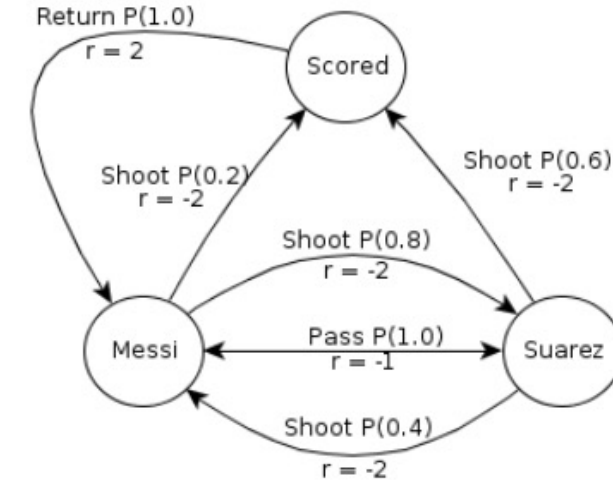|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 | -2 | |
| **V(Suarez)** | 0 | -1 | -1.2 | |
| **V(Scored)** | 0 | 2 | 1 | |

- Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma \, V_i(s') \right]$$

Iteration 2: $V_2(Scored)$

- return



The following diagram shows the transition probabilities and rewards:

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 | -2 | -2.2 |
| **V(Suarez)** | 0 | -1 | -1.2 |  |
| **V(Scored)** | 0 | 2 | 1 |  |



Return P(1.0) r = 2 Scored
Shoot P(0.2) r = -2
Shoot P(0.6) r = -2
Shoot P(0.8) r = -2
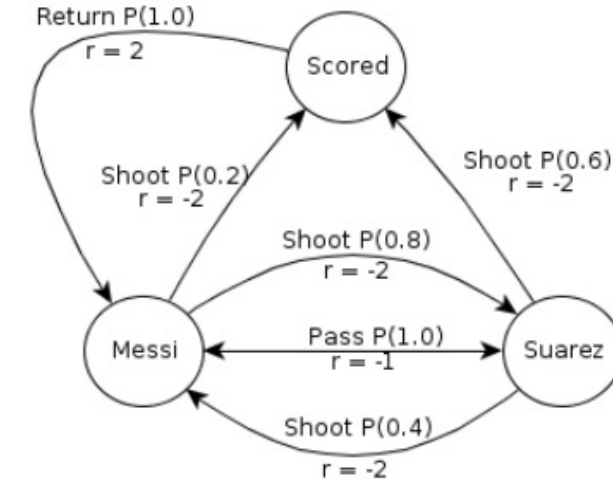Messi
Pass P(1.0) r = -1
Suarez
Shoot P(0.4) r = -2

- Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma V_i(s') \right]$$

$-2.2$

Iteration 3: $V_3(Messi)$

- shoot
$$P_{shoot}(Sua|M)[r(M, shoot, S) + \gamma V(Suarez)]$$
$$+ P_{shoot}(Scored|M)[r(M, shoot, Scored) + \gamma V(Scored)]$$
$$= 0.8[-2 + 1 \times (-1.2)] + 0.2[-2 + 1 \times (1)]$$
$$= -2.76 \quad -2.2$$

- pass
$$P_{pass}(Suarez|M)[r(M, pass, S) + \gamma V(Suarez)] = 1 \times [-1 + 1 \times (-1.2)]$$

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| V(Messi) | 0 | -1 | -2 | -2.2 |
| V(Suarez) | 0 | -1 | -1.2 | -2.2? |
| V(Scored) | 0 | 2 | 1 |  |

- Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma V_i(s') \right]$$

Iteration 3: $V_3(Suarez)$

- shoot

$$P_{shoot}(Messi | Suarez) \left[ r(Suarez, shoot, Messi) + \gamma V(Messi) \right]$$
$$+ P_{shoot}(Scored | Suarez) \left[ r(Suarez, shoot, Scored) + \gamma V(Scored) \right]$$
$$= 0.4 \left[ -2 + 1 \times (-2) \right] + 0.6 \left[ -2 + 1 \times (1) \right] = -2.2$$

max

- pass

$$P_{pass}(Messi | Suarez) \left[ r(Suarez, pass, Messi) + \gamma V(Messi) \right]$$
$$= 1 \times \left[ -1 + 1 \times (-2) \right] = -3$$

Thao Le

18

# Workshop Problem

The following diagram shows the transition probabilities and rewards:

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 | -2 | -2.2 |
| **V(Suarez)** | 0 | -1 | -1.2 | -2.2 |
| **V(Scored)** | 0 | 2 | 1 | 0 |

- Set $V_{i+1}$ to result of Bellman's **right hand side** using $V_i$ in place of $V$:

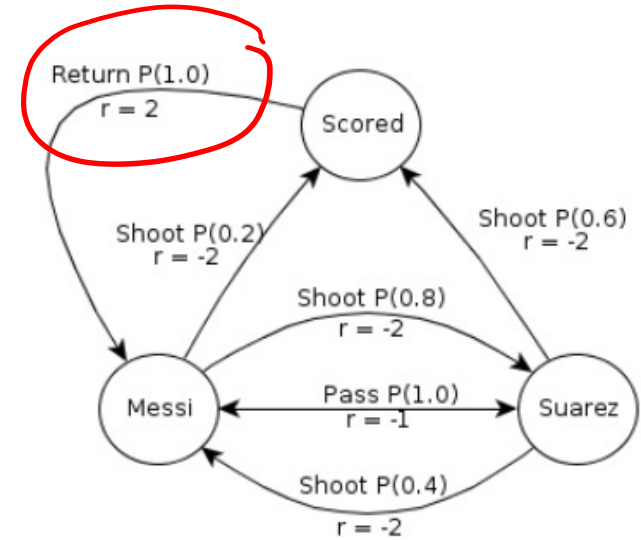$$V_{i+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) \left[ r(s, a, s') + \gamma V_i(s') \right]$$

Iteration 3: $V_3(Scored)$

- return

$$= P_{return}(Messi | Scored) \left[ r(Scored, return, Messi) + \gamma V(Messi) \right]$$

$$= 1 \times \left[ 2 + 1 \times (-2) \right] = 0$$


Return P(1.0) r = 2; Scored; Shoot P(0.2) r = -2; Shoot P(0.6) r = -2; Shoot P(0.8) r = -2; Messi; Pass P(1.0) r = -1; Suarez; Shoot P(0.4) r = -2

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| V(Messi) | 0 | -1 | -2 | -2.2 |
| V(Suarez) | 0 | -1 | -1.2 | -2.2 |
| V(Scored) | 0 | 2 | 1 | 0 |

*from Pass*

*from Shoot*

*from Return*

Return P(1.0)
r = 2

Scored

Shoot P(0.2)
r = -2

Shoot P(0.6)
r = -2

Shoot P(0.8)
r = -2

Messi

Pass P(1.0)
r = -1

Suarez

Shoot P(0.4)
r = -2

If we only have 3 iterations, what actions did we take to maximise the reward?

- Messi Pass
- Suarez Shoot
- Scored Return

# Workshop Problem

|  | Iteration 0 | Iteration 1 | Iteration 2 | Iteration 3 |
|---|---|---|---|---|
| **V(Messi)** | 0 | -1 | -2 | -2.2 |
| **V(Suarez)** | 0 | -1 | -1.2 | -2.2 |
| **V(Scored)** | 0 | 2 | 1 | 0 |

When to stop the iteration?

The iteration is stopped when Δ reaches some pre-defined threshold $\theta$

(when the largest change in the values between iterations is "small enough")



**Algorithm – Value iteration**

**Input:** MDP $M = \langle S, s_0, A, P_a(s' \mid s), r(s, a, s') \rangle$

**Output:** Value function $V$

Set $V$ to arbitrary value function; e.g., $V(s) = 0$ for all $s$

Repeat
  $\Delta \leftarrow 0$
  For each $s \in S$
    $V'(s) \leftarrow \max_{a \in A(s)} \sum_{s' \in S} P_a(s' \mid s) [r(s, a, s') + \gamma V(s')]$
                                  Bellman equation
    $\Delta \leftarrow \max(\Delta, |V'(s) - V(s)|)$
  $V \leftarrow V'$
Until $\Delta \leq \theta$

$\Delta \leq \theta$   $\theta = 0.01$

Messi    Suarez    Scored

① : $\Delta = \max\left(|-1-0|, |-1-0|, |2-0|\right) = 2$

② : $\Delta = \max\left(|-2+1|, |-1.2+1|, |1-2|\right) = 1$

$\Delta \leq \theta$

Thao Le

21