

## Sample Solutions for Problem Set VII: Value Iteration

1. We need to calculate the expected return for each action: pass or shoot.

If Messi passes:

$$\begin{aligned}
 Q(Messi, Pass) &= P_{pass}(Suarez|Messi)[r(Messi, pass, Suarez) + \gamma \cdot V(Suarez)] \\
 &= 1 \cdot [-1 + 1 \cdot -1.2] \\
 &= 1 \cdot -2.2 \\
 &= -2.2
 \end{aligned}$$

If Messi shoots:

$$\begin{aligned}
 Q(Messi, Shoot) &= P_{shoot}(Suarez|Messi)[r(Messi, shoot, Suarez) + \gamma \cdot V(Suarez)] + \\
 &\quad P_{shoot}(Scored|Messi)[r(Messi, shoot, Scored) + \gamma \cdot V(Scored)] \\
 &= 0.8[-2 + 1 \cdot -1.2] + 0.2[-2 + 1 \cdot 1.0] \\
 &= -2.56 + (-0.2) \\
 &= -2.76
 \end{aligned}$$

Therefore, to maximise our reward, Messi should pass.

2. To calculate  $V(Messi)$ , we choose the action that maximises our Q-value (expected future discounted reward):

$$\begin{aligned}
 V(Messi) &= \max(Q(Messi, pass), Q(Messi, shoot)) \\
 &= \max(-2.2, -2.76) \text{ (from previous question)} \\
 &= -2.2
 \end{aligned}$$

For *Scored*, there is only one action, which leads directly to the *Messi* state:

$$\begin{aligned}
 V(Scored) &= P_{return}(Messi|Scored)[r(Scored, return, Messi) + \gamma \cdot V(Messi)] \\
 &= 1[2 + 1 \cdot -2.0] \\
 &= 0
 \end{aligned}$$

For Suarez, the situation is similar to Messi:

$$\begin{aligned}
 V(Suarez) &= \max(Q(Suarez, pass), Q(Suarez, shoot)) \\
 &= \max(P_{pass}(Messi|Suarez)[r(Suarez, pass, Messi) + \gamma \cdot V(Messi), \\
 &\quad (P_{shoot}(Messi|Suarez)[r(Suarez, shoot, Messi) + \gamma \cdot V(Messi) + \\
 &\quad P_{shoot}(Scored|Suarez)[r(Suarez, shoot, Scored) + \gamma \cdot V(Scored)]) \\
 &= \max(1.0[-1 + 1 \cdot -2.0], (0.4[-2 + 1 \cdot 2.0] + 0.6[-2 + 1 \cdot 1.0])) \\
 &= \max(-3, (0.4[-2 + 1 \cdot -2.0] + 0.6[-2 + 1 \cdot 1.0])) \\
 &= \max(-3, (-1.6 + -0.6)) \\
 &= -2.2
 \end{aligned}$$

Thus, the new table is:

Iteration	1	2	3	4
V(Messi)	= 0.0	-1.0	-2.0	-2.2
V(Suarez)	= 0.0	-1.0	-1.2	-2.2
V(Scored)	= 0.0	2.0	1.0	0.0