## The world

A discount factor: $\gamma = 1/2$. The learning rate: $\alpha = 0.1$.

$r(s,a) =$

initial $Q(s,a) =$

$$R = (0,0,100,0,0,0)^T$$
$$Q_0 = (6,4,2,1,3,5)^T$$
$$P_1 = (5,3,1,1,5,6)^T$$

$$Q_1 = Q_0 + 0.1 \cdot (R + 0.5 \cdot P_1 - Q_0) = \begin{pmatrix} 5.65 \\ 3.75 \\ 11.85 \\ 0.95 \\ 2.95 \\ 4.8 \end{pmatrix}$$

$$P_2 = (4.8, 11.85, 0.95, 0.95, 4.8, 5.65)^T$$

$$Q_2 = Q_1 + 0.1 \cdot (R + 0.5 \cdot P_2 - Q_1) = \begin{pmatrix} 5.325 \\ 3.9675 \\ 20.7125 \\ 0.9025 \\ 2.895 \\ 4.6025 \end{pmatrix}$$

$$P_3 = (4.6025, 20.7175, 0.9025, 0.9025, 4.6025, 5.325)^T$$

$$Q_3 = Q_2 + 0.1 \cdot (R + 0.5 \cdot P_3 - Q_2) = \begin{pmatrix} 5.02263 \\ 4.60663 \\ 28.6864 \\ 0.857375 \\ 2.83563 \\ 4.4085 \end{pmatrix}$$

$$\vec{Q}_{\text{true}} =$$

| | 12.5 | | 25 | | 0 |
|---|---|---|---|---|---|
| 25 | | 50 | | 100 | |

$$\vec{R} =$$

| | 0 | | 0 | | 0 |
|---|---|---|---|---|---|
| 0 | | 0 | | 100 | |

$$\vec{Q}_0 =$$

| | 5 | | 3 | | 1 |
|---|---|---|---|---|---|
| 6 | | 4 | | 2 | |

$$\hat{\pi} =$$

(not part of the algorithm)

$$\vec{P}_1 =$$

| | 6 | | 5 | | 1 |
|---|---|---|---|---|---|
| 5 | | 3 | | 1 | |

$$\vec{Q}_1 =$$

| | 4.8 | | 2.95 | | 0.95 |
|---|---|---|---|---|---|
| 5.65 | | 3.74 | | 11.85 | |

$$\hat{\pi} =$$

(not part of the algorithm)

$$\vec{P}_2 =$$

| | 5.65 | | 4.8 | | 0.95 |
|---|---|---|---|---|---|
| 4.8 | | 11.85 | | 0.95 | |

$$\vec{Q}_2 =$$

| | 4.6025 | | 2.895 | | 0.9025 |
|---|---|---|---|---|---|
| 5.325 | | 3.9675 | | 20.7175 | |

$$\hat{\pi} =$$

(not part of the algorithm)

$$\vec{P}_3 =$$

| | 5.325 | | 4.6025 | | 0.9025 |
|---|---|---|---|---|---|
| 4.6025 | | 20.7175 | | 0.9025 | |

$$\vec{Q}_3 =$$

| | 4.4085 | | 2.835625 | | 0.94765 |
|---|---|---|---|---|---|
| 5.022625 | | 4.606625 | | 28.695625 | |

$$\hat{\pi} =$$

(not part of the algorithm)