# Machine Learning

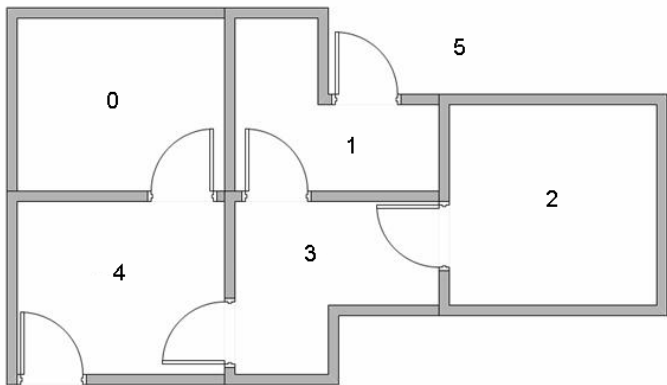## Machine Learning: Jordan Boyd-Graber
University of Maryland
*Q* LEARNING

Slides adapted from John McCullock

**Content Questions**

# Scenario

Goal State

# Scenario: Escape!

# Rewards

**Reward Matrix**

$$
R = \begin{array}{c c}
 & \textbf{Action} \\
\textbf{State} & \begin{array}{c c c c c c}
0 & 1 & 2 & 3 & 4 & 5
\end{array} \\
\begin{array}{c}
0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5
\end{array} &
\begin{bmatrix}
-1 & -1 & -1 & -1 & 0 & -1 \\
-1 & -1 & -1 & 0 & -1 & 100 \\
-1 & -1 & -1 & 0 & -1 & -1 \\
-1 & 0 & 0 & -1 & 0 & -1 \\
0 & -1 & -1 & 0 & -1 & 100 \\
-1 & 0 & -1 & -1 & 0 & 100
\end{bmatrix}
\end{array}
$$

100 Goal

0 Valid Transition

-1 Impossible

## Q-Learning Algorithm

For each $s, a$ initialize table entry $\hat{Q}(s, a) \leftarrow 0$
Observe current state $s$
Do forever:

- Select an action $a$ and execute it
- Receive immediate reward $r$
- Observe the new state $s'$
- Update the table entry for $\hat{Q}(s, a)$ as follows:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$$

- $s \leftarrow s'$

**Initial $Q$ Matrix**

$$Q= \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left[\begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array}\right] \end{array}$$

- Suppose we start in Room 1
- And we'll go to Room 5 afterward

**In Room 5**

$$
\begin{array}{c}
\qquad\qquad\qquad \textbf{Action} \\
\begin{array}{cc}
\textbf{State} & \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\
R = \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} & \begin{bmatrix}
-1 & -1 & -1 & -1 & 0 & -1 \\
-1 & -1 & -1 & 0 & -1 & 100 \\
-1 & -1 & -1 & 0 & -1 & -1 \\
-1 & 0 & 0 & -1 & 0 & -1 \\
0 & -1 & -1 & 0 & -1 & 100 \\
-1 & 0 & -1 & -1 & 0 & 100
\end{bmatrix}
\end{array}
\end{array}
$$

What is the updated $Q$ matrix? ($\gamma = .8$)

$$\hat{Q}(s,a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s',a')$$

**Updated $Q$ for Room $1 \rightarrow 5$**

$$\hat{Q}(1,5) = R(1,5) + \gamma \max\left[\hat{Q}(5,0), \dots \hat{Q}(5,5)\right] \tag{1}$$

**Updated $Q$ for Room $1 \rightarrow 5$**

$$\hat{Q}(1,5) = R(1,5) + \gamma \max\left[\hat{Q}(5,0), \dots \hat{Q}(5,5)\right] \qquad (1)$$

$$\hat{Q}(1,5) = 100 + \gamma \cdot 0 \qquad (2)$$

**Update $Q$ for Room $5 \rightarrow 1$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

$$R = \begin{array}{c} \textbf{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} & & \textbf{Action} & & & \\ 0 & 1 & 2 & 3 & 4 & 5 \\ \left[ \begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array} \right] \end{array}$$

(3)

**Update $Q$ for Room $5 \rightarrow 1$**

$$\hat{Q} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

$$R = \begin{matrix} & \begin{matrix} & & & \textbf{Action} & & \\ \textbf{State} & 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix} \end{matrix}$$

$$\hat{Q}(5,1) = R(5,1) + \gamma \max\big[\hat{Q}(1,0), \dots \hat{Q}(1,5)\big] \tag{3}$$

**Update $Q$ for Room $5 \to 1$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

$$R = \begin{array}{c} \\ \text{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \text{Action} \\ \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left[ \begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array} \right] \end{array}$$

$$\hat{Q}(5,1) = R(5,1) + \gamma \max\left[\hat{Q}(1,0), \ldots \hat{Q}(1,5)\right] \tag{3}$$

$$\hat{Q}(5,1) = 0 + \gamma \cdot 100 \tag{4}$$

**Update $Q$ for Room $1 \to 3$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

$$R = \begin{array}{c} \text{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \text{Action} \\ \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left[ \begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array} \right] \end{array}$$

(5)

**Update $Q$ for Room $1 \rightarrow 3$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left(\begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array}\right) \end{array}$$

Action

| State | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|-----|-----|-----|-----|-----|-----|
| 0 | −1 | −1 | −1 | −1 | 0 | −1 |
| 1 | −1 | −1 | −1 | 0 | −1 | 100 |
| 2 | −1 | −1 | −1 | 0 | −1 | −1 |
| 3 | −1 | 0 | 0 | −1 | 0 | −1 |
| 4 | 0 | −1 | −1 | 0 | −1 | 100 |
| 5 | −1 | 0 | −1 | −1 | 0 | 100 |

$R =$ (for the above matrix)

$$\hat{Q}(1,3) = R(1,3) + \gamma \max\left[\hat{Q}(3,0), \ldots \hat{Q}(3,5)\right] \tag{5}$$

**Update $Q$ for Room $1 \to 3$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left(\begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array}\right) \end{array}$$

| | Action | | | | | |
|---|---|---|---|---|---|---|
| State | 0 | 1 | 2 | 3 | 4 | 5 |
| 0 | −1 | −1 | −1 | −1 | 0 | −1 |
| 1 | −1 | −1 | −1 | 0 | −1 | 100 |
| 2 | −1 | −1 | −1 | 0 | −1 | −1 |
| 3 | −1 | 0 | 0 | −1 | 0 | −1 |
| 4 | 0 | −1 | −1 | 0 | −1 | 100 |
| 5 | −1 | 0 | −1 | −1 | 0 | 100 |

$R=$

$$\hat{Q}(1,3) = R(1,3) + \gamma \max\left[\hat{Q}(3,0), \dots \hat{Q}(3,5)\right] \qquad (5)$$

$$\hat{Q}(1,3) = 0 + \gamma \cdot 0 \qquad (6)$$

**Update $Q$ for Room $3 \rightarrow 4$**

$$\hat{Q} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

$$R = \begin{matrix} & \textbf{Action} \\ \textbf{State} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix} \end{matrix}$$

(7)

**Update $Q$ for Room $3 \rightarrow 4$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left(\begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array}\right) \end{array}$$

Action

| State | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|----|----|----|----|----|----|
| 0 | -1 | -1 | -1 | -1 | 0 | -1 |
| 1 | -1 | -1 | -1 | 0 | -1 | 100 |
| 2 | -1 | -1 | -1 | 0 | -1 | -1 |
| 3 | -1 | 0 | 0 | -1 | 0 | -1 |
| 4 | 0 | -1 | -1 | 0 | -1 | 100 |
| 5 | -1 | 0 | -1 | -1 | 0 | 100 |

$$\hat{Q}(3,4) = R(3,4) + \gamma \max\left[\hat{Q}(4,0), \ldots \hat{Q}(4,5)\right] \tag{7}$$

**Update $Q$ for Room $3 \rightarrow 4$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

Action

| State | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|----|----|----|----|----|----|
| 0 | −1 | −1 | −1 | −1 | 0 | −1 |
| 1 | −1 | −1 | −1 | 0 | −1 | 100 |
| 2 | −1 | −1 | −1 | 0 | −1 | −1 |
| 3 | −1 | 0 | 0 | −1 | 0 | −1 |
| 4 | 0 | −1 | −1 | 0 | −1 | 100 |
| 5 | −1 | 0 | −1 | −1 | 0 | 100 |

$R=$

$$\hat{Q}(3,4) = R(3,4) + \gamma \max \left[ \hat{Q}(4,0), \dots \hat{Q}(4,5) \right] \tag{7}$$

$$\hat{Q}(3,4) = 0 + \gamma \cdot 0 \tag{8}$$

**Update $Q$ for Room $4 \rightarrow 5$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

| | Action | | | | | |
|---|---|---|---|---|---|---|
| **State** | **0** | **1** | **2** | **3** | **4** | **5** |
| 0 | $-1$ | $-1$ | $-1$ | $-1$ | $0$ | $-1$ |
| 1 | $-1$ | $-1$ | $-1$ | $0$ | $-1$ | $100$ |
| $R =$ 2 | $-1$ | $-1$ | $-1$ | $0$ | $-1$ | $-1$ |
| 3 | $-1$ | $0$ | $0$ | $-1$ | $0$ | $-1$ |
| 4 | $0$ | $-1$ | $-1$ | $0$ | $-1$ | $100$ |
| 5 | $-1$ | $0$ | $-1$ | $-1$ | $0$ | $100$ |

(9)

**Update $Q$ for Room $4 \rightarrow 5$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{pmatrix} \end{array}$$

$$R = \begin{array}{c} \textbf{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} & & \textbf{Action} & & & \\ 0 & 1 & 2 & 3 & 4 & 5 \\ \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix} \end{array}$$

$$\hat{Q}(4,5) = R(4,5) + \gamma \max \left[ \hat{Q}(5,0), \ldots \hat{Q}(5,5) \right] \tag{9}$$

**Update $Q$ for Room $4 \rightarrow 5$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{pmatrix} \end{array}$$

$$R = \begin{array}{c} \text{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \text{Action} \\ \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix} \end{array}$$

$$\hat{Q}(4,5) = R(4,5) + \gamma \max\left[\hat{Q}(5,0), \ldots \hat{Q}(5,5)\right] \tag{9}$$

$$\hat{Q}(4,5) = 100 + \gamma \cdot 80 \tag{10}$$

**Update $Q$ for Room $5 \rightarrow 4$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 164 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

$$R = \begin{array}{c} \textbf{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} \textbf{Action} \\ 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left[ \begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array} \right] \end{array}$$

(11)

**Update $Q$ for Room $5 \rightarrow 4$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 164 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

Action

| State | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|----|----|----|----|----|-----|
| 0 | -1 | -1 | -1 | -1 | 0 | -1 |
| 1 | -1 | -1 | -1 | 0 | -1 | 100 |
| 2 | -1 | -1 | -1 | 0 | -1 | -1 |
| 3 | -1 | 0 | 0 | -1 | 0 | -1 |
| 4 | 0 | -1 | -1 | 0 | -1 | 100 |
| 5 | -1 | 0 | -1 | -1 | 0 | 100 |

$$\hat{Q}(5,4) = R(5,4) + \gamma \max \left[ \hat{Q}(4,0), \ldots \hat{Q}(4,5) \right] \tag{11}$$

**Update $Q$ for Room $5 \to 4$**

$$\hat{Q} = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left( \begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 164 \\ 0 & 80 & 0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

$$R = \begin{array}{c} \\ \text{State} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{c} \text{Action} \\ \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \end{array} \\ \left[ \begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array} \right] \end{array}$$
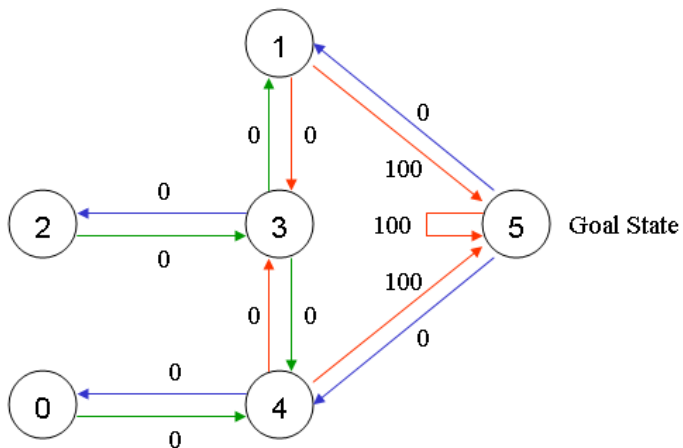
$$\hat{Q}(5,4) = R(5,4) + \gamma \max \left[ \hat{Q}(4,0), \dots \hat{Q}(4,5) \right] \tag{11}$$

$$\hat{Q}(5,4) = 0 + \gamma \cdot 164 = 131 \tag{12}$$

**If you keep going . . .**

$$
Q = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array}
\begin{array}{cccccc}
0 & 1 & 2 & 3 & 4 & 5 \\
\begin{bmatrix}
0 & 0 & 0 & 0 & 400 & 0 \\
0 & 0 & 0 & 320 & 0 & 500 \\
0 & 0 & 0 & 320 & 0 & 0 \\
0 & 400 & 256 & 0 & 400 & 0 \\
320 & 0 & 0 & 320 & 0 & 500 \\
0 & 400 & 0 & 0 & 400 & 500
\end{bmatrix}
\end{array}
$$

**Is this really hard?**

- Creating the state space
- Estimating rewards
- Choosing action with incomplete learning