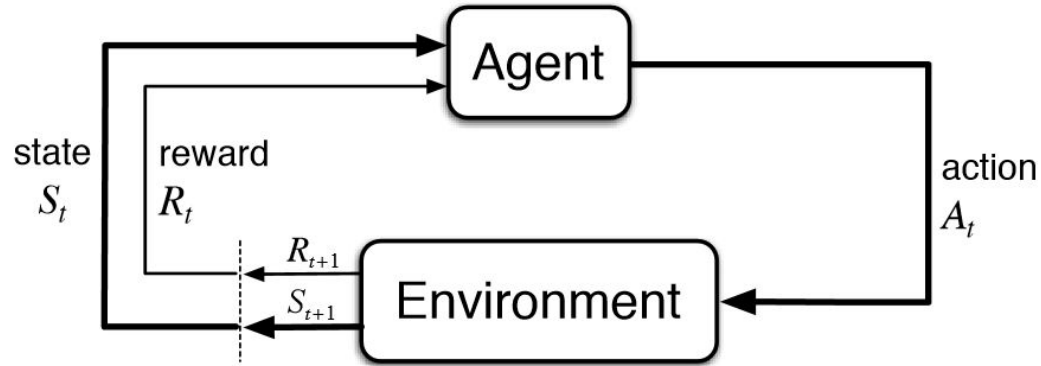


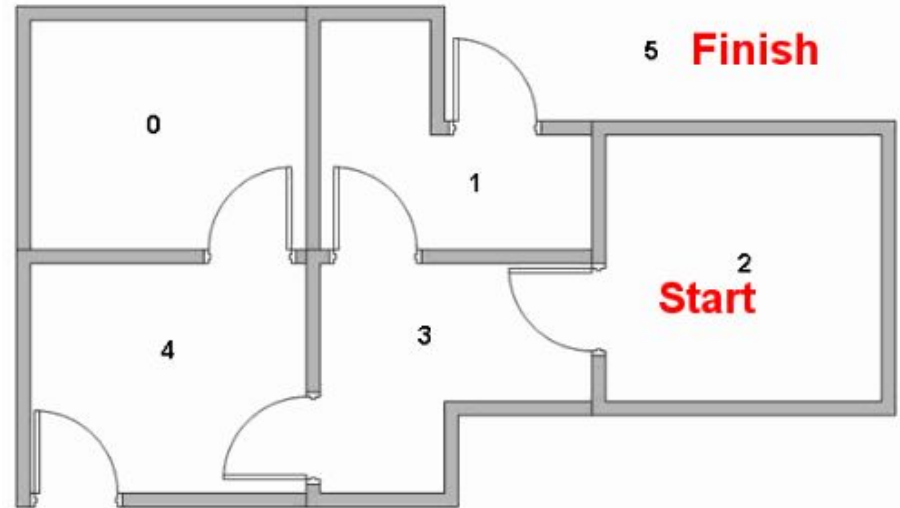
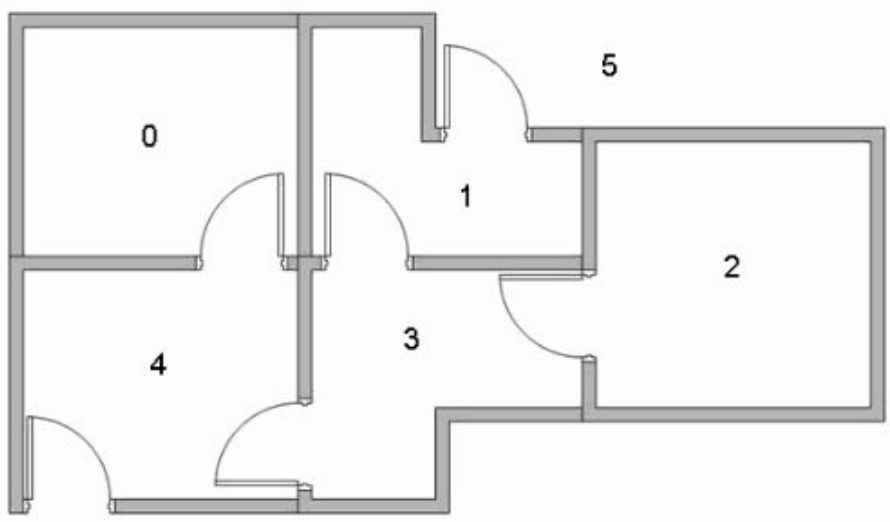
Reinforcement Learning

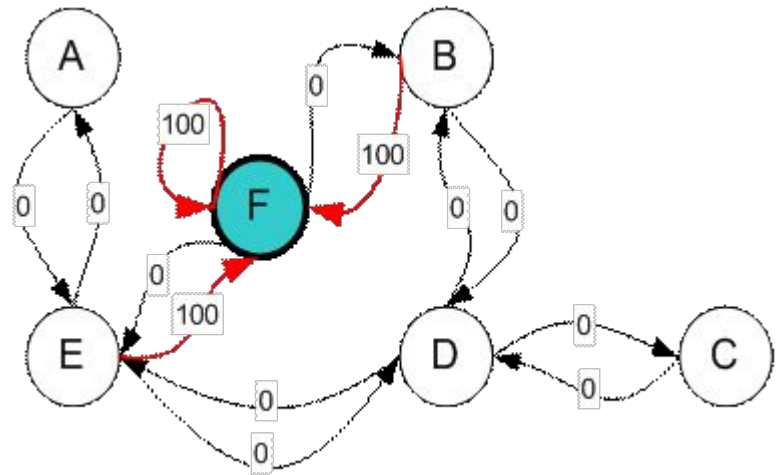
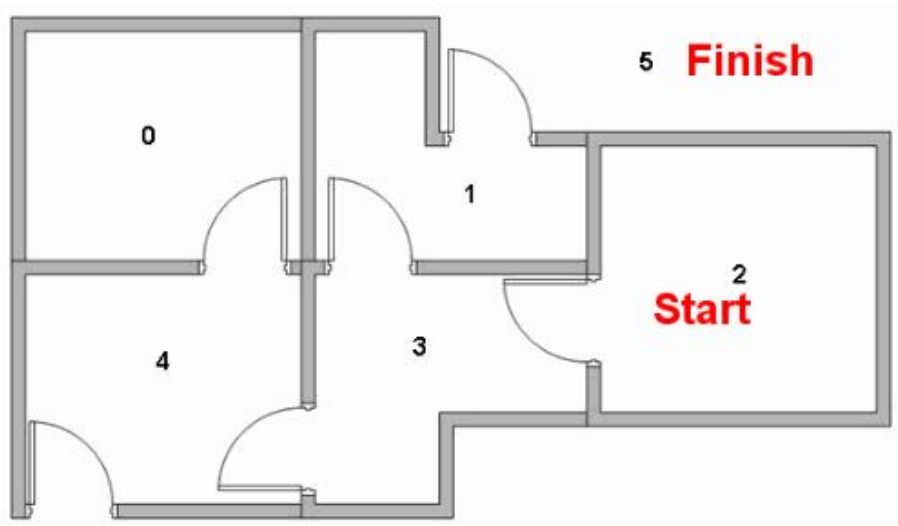
Q-Learning

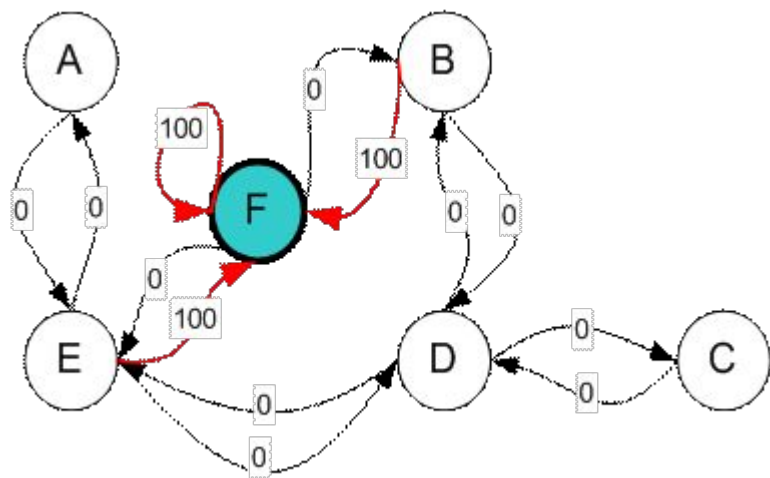
Reinforcement Learning



Numerical example







$\mathbf{R} =$

<i>state \ action</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>A</i>	–	–	–	–	0	–
<i>B</i>	–	–	–	0	–	100
<i>C</i>	–	–	–	0	–	–
<i>D</i>	–	0	0	–	0	–
<i>E</i>	0	–	–	0	–	100
<i>F</i>	–	0	–	–	0	100

Q Learning

1. Set the gamma parameter and environment rewards in matrix R.
2. Initialize matrix Q to zero.
3. For each episode:

 Select a random initial state.

 Do While the goal state hasn't been reached.

 Select one among all possible actions for the current state.

 Using this possible action, consider going to the next state.

 Get maximum Q value for this next state based on all possible actions.

 Compute: $Q(\text{state}, \text{action}) = R(\text{state}, \text{action}) + \text{Gamma} * \text{Max}[Q(\text{next state}, \text{all actions})]$

 Set the next state as the current state.

 End Do

End For

Episode 1

Lets us set the value of learning $\gamma=0.8$

Lets suppose that we start in state B

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix} \quad \mathbf{R} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \end{matrix}$$

$$\mathbf{Q}(\text{state}, \text{action}) = \mathbf{R}(\text{state}, \text{action}) + \gamma \cdot \text{Max}[\mathbf{Q}(\text{next state}, \text{all actions})]$$

$$\mathbf{Q}(B, F) = \mathbf{R}(B, F) + 0.8 \cdot \text{Max}\{\mathbf{Q}(F, B), \mathbf{Q}(F, E), \mathbf{Q}(F, F)\} = 100 + 0.8 \cdot 0 = 100$$

Episode 1

$$\mathbf{Q}(\text{state}, \text{action}) = \mathbf{R}(\text{state}, \text{action}) + \gamma \cdot \text{Max}[\mathbf{Q}(\text{next state}, \text{all actions})]$$

$$\mathbf{Q}(B, F) = \mathbf{R}(B, F) + 0.8 \cdot \text{Max}\{\mathbf{Q}(F, B), \mathbf{Q}(F, E), \mathbf{Q}(F, F)\} = 100 + 0.8 \cdot 0 = 100$$

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

F is now the current state. Because F is the goal state, we finish one episode

Episode 2

Lets suppose that we start in state D

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

$$\mathbf{R} = \begin{matrix} & \begin{matrix} state \backslash action & A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \end{matrix}$$

$$\mathbf{Q}(state, action) = \mathbf{R}(state, action) + \gamma \cdot \text{Max}[\mathbf{Q}(next\ state, all\ actions)]$$

$$\mathbf{Q}(D, B) = \mathbf{R}(D, B) + 0.8 \cdot \text{Max}\{\mathbf{Q}(B, D), \mathbf{Q}(B, F)\} = 0 + 0.8 \cdot \text{Max}\{0, 100\} = 80$$

Episode 2

Lets suppose that we start in state D

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

$$\mathbf{R} = \begin{matrix} & \begin{matrix} state \backslash action & A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} - & - & - & - & 0 & - \\ - & - & - & 0 & - & 100 \\ - & - & - & 0 & - & - \\ - & 0 & 0 & - & 0 & - \\ 0 & - & - & 0 & - & 100 \\ - & 0 & - & - & 0 & 100 \end{bmatrix} \end{matrix}$$

$$\mathbf{Q}(state, action) = \mathbf{R}(state, action) + \gamma \cdot \text{Max}[\mathbf{Q}(next\ state, all\ actions)]$$

$$\mathbf{Q}(D, B) = \mathbf{R}(D, B) + 0.8 \cdot \text{Max}\{\mathbf{Q}(B, D), \mathbf{Q}(B, F)\} = 0 + 0.8 \cdot \text{Max}\{0, 100\} = 80$$

Episode 2

$$\mathbf{Q}(\textit{state}, \textit{action}) = \mathbf{R}(\textit{state}, \textit{action}) + \gamma \cdot \textit{Max}[\mathbf{Q}(\textit{next state}, \textit{all actions})]$$

$$\mathbf{Q}(D, B) = \mathbf{R}(D, B) + 0.8 \cdot \textit{Max}\{\mathbf{Q}(B, D), \mathbf{Q}(B, F)\} = 0 + 0.8 \cdot \textit{Max}\{0, 100\} = 80$$

$$\mathbf{Q} = \begin{array}{c} \begin{array}{cccccc} & A & B & C & D & E & F \end{array} \\ \begin{array}{l} A \\ B \\ C \\ D \\ E \\ F \end{array} \end{array} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Episode 2

The next state B is now the current state

$$\mathbf{Q} = \begin{array}{c|cccccc} & A & B & C & D & E & F \\ \hline A & 0 & 0 & 0 & 0 & 0 & 0 \\ B & 0 & 0 & 0 & 0 & 0 & 100 \\ C & 0 & 0 & 0 & 0 & 0 & 0 \\ D & 0 & 80 & 0 & 0 & 0 & 0 \\ E & 0 & 0 & 0 & 0 & 0 & 0 \\ F & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \quad \mathbf{R} = \begin{array}{c|cccccc} \text{state} \backslash \text{action} & A & B & C & D & E & F \\ \hline A & - & - & - & - & 0 & - \\ B & - & - & - & 0 & - & 100 \\ C & - & - & - & 0 & - & - \\ D & - & 0 & 0 & - & 0 & - \\ E & 0 & - & - & 0 & - & 100 \\ F & - & 0 & - & - & 0 & 100 \end{array}$$

$$\mathbf{Q}(\text{state}, \text{action}) = \mathbf{R}(\text{state}, \text{action}) + \gamma \cdot \text{Max}[\mathbf{Q}(\text{next state}, \text{all actions})]$$

$$\begin{aligned} \mathbf{Q}(B, F) &= \mathbf{R}(B, F) + 0.8 \cdot \text{Max}\{\mathbf{Q}(F, B), \mathbf{Q}(F, E), \mathbf{Q}(F, F)\} \\ &= 100 + 0.8 \cdot \text{Max}\{0, 0, 0\} = 100 \end{aligned}$$

Episode 2

$$\mathbf{Q}(\text{state}, \text{action}) = \mathbf{R}(\text{state}, \text{action}) + \gamma \cdot \text{Max}[\mathbf{Q}(\text{next state}, \text{all actions})]$$

$$\mathbf{Q}(B, F) = \mathbf{R}(B, F) + 0.8 \cdot \text{Max}\{\mathbf{Q}(F, B), \mathbf{Q}(F, E), \mathbf{Q}(F, F)\} = 100 + 0.8 \cdot 0 = 100$$

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

F is now the current state. Because F is the goal state, we finish one episode

Convergence

If our agent learn more and more experience through many episodes, it will finally reach convergence values of Q matrix as

- This Q matrix, then can be normalized into %, dividing all valid entries with the highest number

Q =

<i>state \ action</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>A</i>	-	-	-	-	400	-
<i>B</i>	-	-	-	320	-	500
<i>C</i>	-	-	-	320	-	-
<i>D</i>	-	400	256	-	400	-
<i>E</i>	320	-	-	320	-	500
<i>F</i>	-	400	-	-	400	500

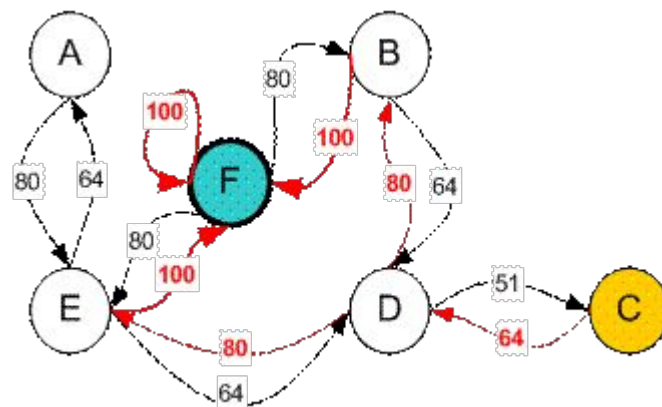
\hat{Q} =

<i>state \ action</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>A</i>	-	-	-	-	80	-
<i>B</i>	-	-	-	64	-	100
<i>C</i>	-	-	-	64	-	-
<i>D</i>	-	80	51	-	80	-
<i>E</i>	64	-	-	64	-	100
<i>F</i>	-	80	-	-	80	100

Convergence

$\hat{Q} =$

<i>state \ action</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>A</i>	-	-	-	-	80	-
<i>B</i>	-	-	-	64	-	100
<i>C</i>	-	-	-	64	-	-
<i>D</i>	-	80	51	-	80	-
<i>E</i>	64	-	-	64	-	100
<i>F</i>	-	80	-	-	80	100












Robot Game

In the case of the robot game, to reiterate the scoring/reward structure is:

- power = +1
- mine = -100
- end = +100

$$\text{New } Q(s, a) = Q(s, a) + \alpha [R(s, a) + \gamma \max_{a'} Q'(s', a') - Q(s, a)]$$

- New Q Value for that state and the action
- Learning Rate
- Reward for taking that action at that state
- Current Q Values
- Maximum expected future reward given the new state (s') and all possible actions at that new state.
- Discount Rate

					
					
					
					
				End	

Actions : ↑ → ↓ ←

Start	0	0	0	0
Nothing / Blank	0	0	0	0
Power	0	0	0	0
Mines	0	0	0	0
END	0	0	0	0

New $Q(\text{start}, \text{right}) = Q(\text{start}, \text{right}) + \alpha [\text{some ... Delta value}]$

Some ... Delta value = $R(\text{start}, \text{right}) + \max\{Q'(\text{nothing}, \text{down}), Q'(\text{nothing}, \text{left}), Q'(\text{nothing}, \text{right})\} - Q(\text{start}, \text{right})$

Some ... Delta value = $0 + 0.9 * 0 - 0 = 0$

New $Q(\text{start}, \text{right}) = 0 + 0.1 * 0 = 0$