



(<http://people.revoledu.com/kardi/>)

MENU

Q-Learning Numerical Example

by [Kardi Teknomo \(../index.html\)](#)



([purchase.html](#))

< [Previous \(Q-Learning-Algorithm.htm\)](#) | [Next \(Tower-of-Hanoi.htm\)](#) | [Contents \(index.html\)](#) >

[Read this tutorial comfortably off-line. Click here to purchase the complete E-book of this tutorial \(purchase.html\)](#)

Q-Learning Numerical Examples

To understand how the [Q \(Q-Learning-Algorithm.htm\) learning algorithm \(Q-Learning-Algorithm.htm\)](#) works, we will go through several steps of numerical examples. The rest of the steps can be can be confirm using the program that I made ([the companion files of this tutorial can be purchasecan be purchase if you purchase here \(purchase.html\)](#))

Let us set the value of learning parameter $\gamma = 0.8$ and initial state as room B.

First we set matrix Q as a zero matrix.

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

I put again the instant reward matrix R that represents the environment in here for your convenience.

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

Look at the second row (state B) of matrix R. There are two possible actions for the current state B, that is to go to state D, or go to

state F. By random selection, we select to go to F as our action.

Now we consider that suppose we are in state F. Look at the sixth row of reward matrix R (i.e. state F). It has 3 possible actions to go to state B, E or F.

$$Q(state, action) = R(state, action) + \gamma \cdot Max[Q(next\ state, all\ actions)]$$
$$Q(B, F) = R(B, F) + 0.8 \cdot Max\{Q(F, B), Q(F, E), Q(F, F)\} = 100 + 0.8 \cdot 0 = 100$$

Since matrix Q that is still zero, $Q(F, B), Q(F, E), Q(F, F)$ are all zero. The result of computation $Q(B, F)$ is also 100 because of the instant reward.

The next state is F, now become the current state. Because F is the goal state, we finish one episode. Our agent's brain now contain updated matrix Q as

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

For the next episode, we start with initial random state. This time for instance we have state D as our initial state.

Look at the fourth row of matrix R; it has 3 possible actions, that is to go to state B, C and E. By random selection, we select to go to state B as our action.

Now we imagine that we are in state B. Look at the second row of reward matrix R (i.e. state B). It has 2 possible actions to go to state D or state F. Then, we compute the Q value

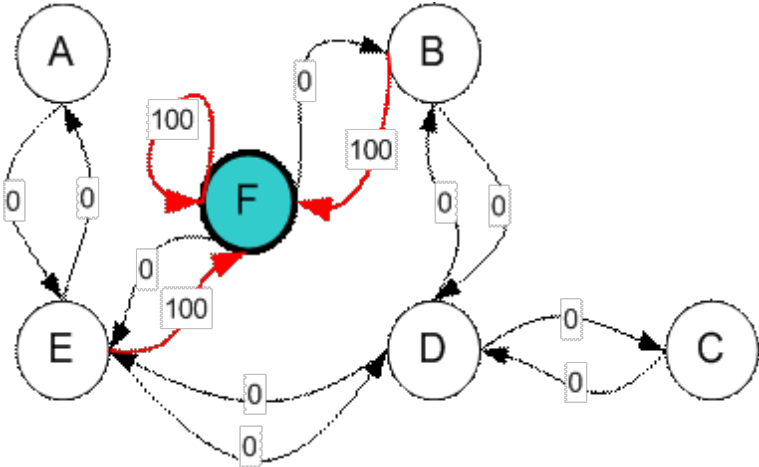
$$Q(state, action) = R(state, action) + \gamma \cdot Max[Q(next\ state, all\ actions)]$$
$$Q(D, B) = R(D, B) + 0.8 \cdot Max\{Q(B, D), Q(B, F)\} = 0 + 0.8 \cdot Max\{0, 100\} = 80$$

We use the updated matrix Q from the last episode. $Q(B, D) = 0$ and $Q(B, F) = 100$. The result of computation $Q(D, B) = 80$ because of the reward is zero. The Q matrix becomes

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

The next state is B, now become the current state. We repeat the inner loop in Q learning algorithm because state B is not the goal state.

For the new loop, the current state is state B. I copy again the state diagram that represent instant reward matrix R for your convenient.



There are two possible actions from the current state B, that is to go to state D, or go to state F. By lucky draw, our action selected is state F.

Q-Learning: Now we think of state F that has 3 possible actions to go to state B, E or F. We compute the Q value using the maximum value of these possible actions.

$$Q(state, action) = R(state, action) + \gamma \cdot Max[Q(next\ state, all\ actions)]$$
$$Q(B, F) = R(B, F) + 0.8 \cdot Max\{Q(F, B), Q(F, E), Q(F, F)\}$$
$$= 100 + 0.8 \cdot Max\{0, 0, 0\} = 100$$

The entries of updated Q matrix contain $Q(F, B), Q(F, E), Q(F, F)$ are all zero. The result of computation $Q(B, F)$ is also 100 because of the instant reward. This result does not change the Q matrix.

Because F is the goal state, we finish this episode. Our agent's brain now contain updated matrix Q as

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

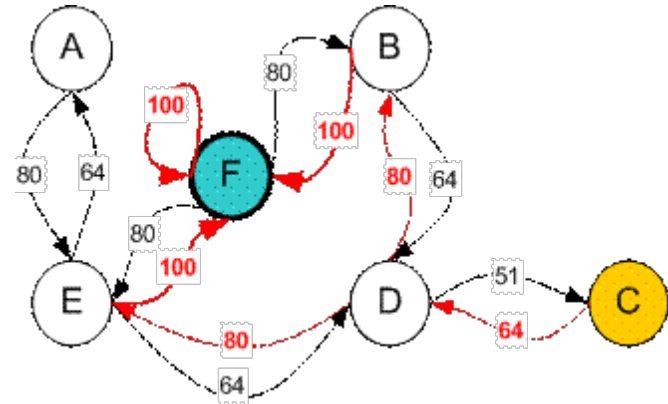
If our agent learns more and more experience through many episodes, it will finally reach convergence values of Q matrix as

$$Q = \begin{matrix} state \backslash action & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} - & - & - & - & 400 & - \\ - & - & - & 320 & - & 500 \\ - & - & - & 320 & - & - \\ - & 400 & 256 & - & 400 & - \\ 320 & - & - & 320 & - & 500 \\ - & 400 & - & - & 400 & 500 \end{bmatrix} \end{matrix}$$

This Q matrix, then can be normalized into a percentage by dividing all valid entries with the highest number (divided by 500 in this case) becomes

$$\hat{Q} = \begin{matrix} state \backslash action & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} - & - & - & - & 80 & - \\ - & - & - & 64 & - & 100 \\ - & - & - & 64 & - & - \\ - & 80 & 51 & - & 80 & - \\ 64 & - & - & 64 & - & 100 \\ - & 80 & - & - & 80 & 100 \end{bmatrix} \end{matrix}$$

Once the Q matrix reaches almost the convergence value, our agent can reach the goal in an optimum way. To trace the sequence of states, it can easily compute by finding action that makes maximum Q for this state.



For example from initial State C, it can use the Q matrix as follow:

From State C the maximum Q produces action to go to state D

From State D the maximum Q has two alternatives to go to state B or E. Suppose we choose arbitrary to go to B

From State B the maximum value produces action to go to state F

Thus the sequence is C -> D -> B -> F

< [Previous \(Q-Learning-Algorithm.htm\)](#) | [Next \(Tower-of-Hanoi.htm\)](#) | [Contents \(index.html\)](#) >

[This tutorial is copyrighted. \(../..copyright.html\)](#)

Preferable reference for this tutorial is

Teknomo, Kardi. 2005. Q-Learning by Examples. <http://people.revoledu.com/kardi/tutorial/ReinforcementLearning/index.html>

Copyright © 2017 Kardi Teknomo
Revoledu Design