

Reinforced Machine Learning using Q-Algorithm Documentation

Author: Stiliyan Tonev
Contact: stiliyan21[at]gmail.com

Content:

- Abstract

1. What is Q-Learning?
2. How can it be used?
3. Learning how it works.
4. My Implementation & how to make it run.
5. What's next?

Abstract

Reinforced learning is a type of intelligent computer agent, which has the ability to observe and adapt to the environment he is put in. Here is an implementation of the Q-Algorithm which is used for such agents. It is not the only one but it is a stepping stone for other, rival algorithms. This implementation was built with generality in mind, so the developer can make multi-dimensional environment without too much effort.

What is Q-Learning?

Q-Learning is a type of reinforced learning algorithm. Reinforced means it observes and adapts to the environment. It is used in Artificial Intelligence (AI) systems in order to simulate human like process of learning. It is related with the so-called “Markov Chains”. They represent the math side of the algorithm, but in many places it is described by logic and is quite intuitive, which makes it easier and more understandable for beginners.

How can it be used?

The algorithm finds application where the “agent” is left unsupervised and is expected to reach a given conclusion. For example a robot, of those that clean the floor, during the process of cleaning it can bump in to an obstacle, or reach a point from which it can't return (stairs), so by using the Q-Algorithm the robot will learn that those places are dangerous and won't pass through them anymore.

Another example can be creating a picture by a given model, the “agent” will learn about the given model in the first few tries to recreate it, collecting information about the model's general properties, and later use that information to generate something which follows those general properties.

It is also used in order to find paths, as with every new cycle it learns which path is better than the previous, and after learning enough it will be able to determine the optimal path to a given objective.

Learning how it works.

The structure of the algorithm is as follows:

- $S \rightarrow$ States in which the “agent” can go, those who are forbidden too, each state has a “reward”.
- $A \rightarrow$ Actions which the “agent” can take, each state has a set of actions the agent can take, and each action has a value, and an exit state.
- Learning coefficient (α) $[0,1] \rightarrow$ It plays a role like a believe factor, if it is 1 or larger the “agent” will believe he is on the right way evin if he is not.
- Discount Factor (γ) $[0,1] \rightarrow$ Tells the agent how much will the next action matter, it can also lead to the wrong path if too high.
- $s \rightarrow$ Current state of the “agent”.
- $a \rightarrow$ Action choosed from the current state.
- s' (next state) \rightarrow The state in which is the “agent” after executing action “a” in state “s”.
- a' (next action) \rightarrow The choosed action in state s' .

And the o-so-holy Q-Function:

- $Q(s,a) = r + \alpha * \gamma * (\max_{a'}(s'))$
- Which translated means “the Quality of action ‘a’ from state ‘s’ is equal to the (r)Reward in state s’ plus (learning coeficent)*(discount)*(the action with most value from s’)”
- Simply said follow the bigger value.

Now as we have “undurstood” the math lets look at how it works with the enviroment.

We will work it out on a 3x3 grid. Don’t worry if you are confused it will get clearer the more you read it.

0.1111	1.0	-1.0
0.5	0.0	-1.0
2.0	0.3	0.0

$\alpha=0.8$ $\gamma=0.5$

Now our agent start from the top-left cell (0.1111). It chooses the action by comparing the actions value, for our conviniance we will use the reward of the next state as action value so,

Right: 1.0 Down:0.5

It chooses the larger one (Right), and updates the transition value
 $[Right] = Q(\text{top-left}, \text{right}) = 1.0(\text{reward}) + 0.8 * 0.5 * (\max(0.0, 0.1111, -1.0)) = 1.0 + 0.8 * 0.5 * 0 = 1 + 0.40 = 1.40$

And now from the top-center cell:

Down: 0.0 Right:-1.0

Left:(if you go back there you enter endless cycle)

Now the largest is 0.0:

$[Down] = Q(\text{top-center}, \text{down}) = 0.0 + 0.8 * 0.5 * (\max(0.5, 0.3, 0.0)) = 1.0 + 0.8 * 0.5 * 0.5 = 1 + 0.20 = 1.20$

And from there it must be clear that we will end up in the bottom-left cell, because we will move to the biggest neighbour (0.5) and the next move with most value is (2.0).

Indeed this is a very rought example, but the thing you must remember/understand:

{choose action, estimate $Q(s,a)$, assign to the present action, move to s' , repeat}

As we said before that Q in Q-Learning is for Quality which is estimated by own, and neighbour value.

My implementation and how to make it run

My implementation is a java package, containing 3 main files.

1. Mind
2. States
3. Actions

The file Actions.java is the shortest and most simple one.

It contains the declaration of “Actions” their property “Reward” (Value), Name, and EndState.

The value is the “Reward” you get after executing this action.

The end state can be different or remain the same.

And the name is used to make it simpler and understandable.

(Action:

 Name=”Go to university”

 Reward=0.001

 EndState=”University”

Action:

 Name=”Get coffee”

 Reward=2.999

 EndState=BegginState

)

It can be constructed without arguments and using the setters assing values to, or called with its different constructors.

Properties:

- String Name;

- double Reward;
- States EndState;

Constructors:

- Actions();
- Actions(EndState);
- Actions(name,endstate);

Methods:

- void setName("Other action name");
- String getName()
- double getReward();
- void setReward(1.0);
- void setState(endstate);
- States getState();

States.java is a little bit more complex file (but wait till we get to mind).

Every state has "Attributes" which by type is a HashMap → [String,Double]

Properties:

- List<Actions> Actions; → Actions that can be taken from this state.
- HashMap<String,Double> Attributes; → Allows the developer to set different (custom) attributes to actions. Like you want a "checked" marker to keep track of checked cells just add it in to the HashMap and use it when printing. Use your imagination.

Constructors:

- States();
- States(List<String> str,List<Double> arg); → Adds them to the HashMap as each string corresponds to the same index in arg.

Methods:

- Actions `getBestAction()`; → Returns the action with most value.
- Actions `getImprovedAction()`; → Same as `getBeastAction()` except it omits actions that will lead to unsatisfying states no matter they're value.
- Actions `getRandomAction()`; → Used for random exploration.
- `void setAttribute("attribute name",value)` → Sets the value of the given attribute to the given.
- `double getAttribute("attribute name")` → Returns the value of the attribute with the given name.
- `void setAction(Actions)`; → Adds a new action to the List, if there is already an Action with the same name, it will be replaced with the present one.
- `void ForceSuccess()`; → Removes the actions leading to unwanted state, this function does not affect random exploration.
- `void setGoal()` → Sets the "exit" attribute to 1, therefore this is an exit state and is satisfying one.
- `void setObstacle()` → Sets the "exit" attribute to -1, therefore an exit state but not a satisfying one.

Mind.java is the file that puts all of these together.

It has a few functions to help testing or rather, are a Proof-Of-Concept. For more complex models you would have to write the environment by yourself. The default supported environment is an integer matrix, with 0 for empty cell, -1 for blocked/malicious cell, and 1 for goal.

Properties:

- `List<States> DStates;`
- `HashMap<String,Double> Attributes;`

Methods:

- `void setAttribute("name",value);` → Assigns the current value to the given name.
- `double Use("name")` → Return the value of the attribute with the given name.
- `void addState()` → Adds a new state to DStates.
- `double Bellman(s,a);` → Returns the $Q(s,a)$ value by applying the Bellman Equation.
- `void setTable(int[][] arr)` → Accepts an integer matrix with 0 for empty cell, -1 for blocked cell, and 1 for goal.
- `void setRandomTable(s,d,n)` → Generates a random environment of size $(s*s)$, with obstacle density (d) which is advised for you to first see how it works with values in $[0.1,0.3]$ and then choose your preferred, and (n) goals set randomly in the environment.
- `void GraphicSolution();` → Prints the chosen path from a random point to a goal and represents it using Swing library.