



Quality and Innovation

PERFORMANCE EXCELLENCE & DATA SCIENCE FOR DIGITAL TRANSFORMATION IN INDUSTRY 4.0

SATURDAY, MARCH 24TH, 2018

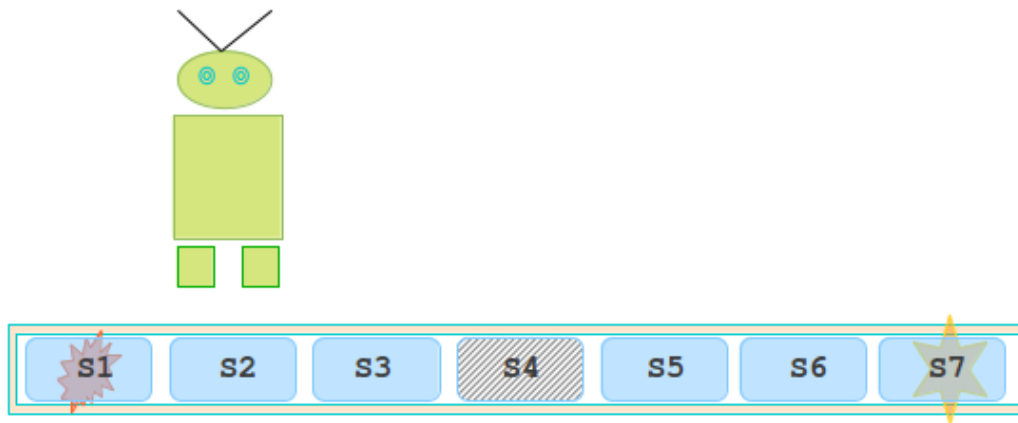
DATA SCIENCE

Reinforcement Learning: Q-Learning with the Hopping Robot

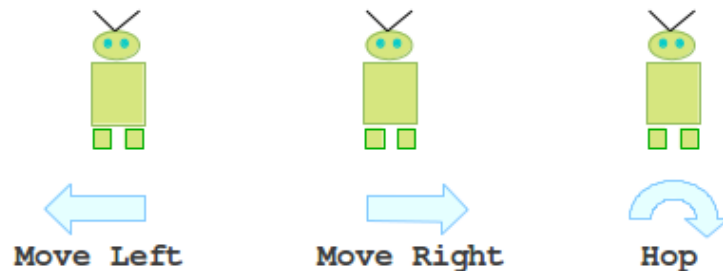
BY NICOLE RADZIWIŁŁ ON DECEMBER 24, 2017 • (5 COMMENTS)

Overview: Reinforcement learning uses “reward” signals to determine how to navigate through a system in the most valuable way. (I’m particularly interested in the variant of reinforcement learning called “Q-Learning” because the goal is to create a “Quality Matrix” that can help you make the best sequence of decisions!) I found a toy robot navigation problem on the web that was solved using custom R code for reinforcement learning, and I wanted to reproduce the solution in different ways than the original author did. This post describes different ways that I solved the problem described at <http://bayesianthink.blogspot.com/2014/05/hopping-robots-and-reinforcement.html> (<http://bayesianthink.blogspot.com/2014/05/hopping-robots-and-reinforcement.html>)

The Problem: Our *agent*, the robot, is placed at random on a board of wood. There’s a hole at s1, a sticky patch at s4, and the robot is trying to make appropriate decisions to navigate to s7 (the target). The image comes from the blog post linked above.



Possible Robot Actions



To solve a problem like this, you can use MODEL-BASED approaches if you know how likely it is that the robot will move from one state to another (that is, the *transition probabilities* for each action) or MODEL-FREE approaches (you don't know how likely it is that the robot will move from state to state, but you *can* figure out a reward structure).

- **Markov Decision Process (MDP)** – If you know the states, actions, rewards, and transition probabilities (which are probably different for each action), you can determine the optimal policy or “path” through the system, given different starting states. (If transition probabilities have nothing to do with decisions that an agent makes, your MDP reduces to a **Markov Chain**.)
- **Reinforcement Learning (RL)** – If you know the states, actions, and rewards (but not the transition probabilities), you can still take an unsupervised approach. Just randomly create lots of hops through your system, and use them to update a matrix that describes the average value of each hop within the context of the system.

Solving a RL problem involves finding the optimal value functions (e.g. the Q matrix in Attempt 1) or the optimal policy (the State-Action matrix in Attempt 2). Although there are many techniques for reinforcement learning, we will use Q-learning because we don't know the transition probabilities for each action. (If we did, we'd model it as a Markov Decision Process and use the MDPtoolbox package instead.) Q-Learning relies on traversing the system in many ways to update a matrix of average expected rewards from each state transition. This equation that it uses is from <https://www.is.uni-freiburg.de/ressourcen/business->

analytics/13_reinforcementlearning.pdf (https://www.is.uni-freiburg.de/ressourcen/business-analytics/13_reinforcementlearning.pdf):

$$Q(s,a) \leftarrow \underbrace{Q(s,a)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \left[\underbrace{r'}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \underbrace{\max_{a'} Q(s',a')}_{\text{expected optimal value}} - \underbrace{Q(s,a)}_{\text{old value}} \right]$$

For this to work, *all states* have to be visited a sufficient number of times, and *all state-action pairs* have to be included in your experience sample. So keep this in mind when you're trying to figure out how many iterations you need.

Attempt 1: Quick Q-Learning with qlearn.R

- **Input:** A rewards matrix R. (That's all you need! Your states are encoded in the matrix.)
- **Output:** A Q matrix from which you can extract optimal policies (or paths) to help you navigate the environment.
- **Pros:** Quick and *very easy*. **Cons:** Does not let you set epsilon (% of random actions), so all episodes are determined randomly and it may take longer to find a solution. Can take a long time to converge.

Set up the rewards matrix so it is a square matrix with *all the states down the rows*, starting with the first and *all the states along the columns*, starting with the first:

```
1 | hopper.rewards <- c(-10, 0.01, 0.01, -1, -1, -1, -1,
2 |   -10, -1, 0.1, -3, -1, -1, -1,
3 |   -1, 0.01, -1, -3, 0.01, -1, -1,
4 |   -1, -1, 0.01, -1, 0.01, 0.01, -1,
5 |   -1, -1, -1, -3, -1, 0.01, 100,
6 |   -1, -1, -1, -1, 0.01, -1, 100,
7 |   -1, -1, -1, -1, -1, 0.01, 100)
8 |
9 | HOP <- matrix(hopper.rewards, nrow=7, ncol=7, byrow=TRUE)
10 | > HOP
11 |      [,1] [,2] [,3] [,4] [,5] [,6] [,7]
12 | [1,] -10  0.01 0.01 -1 -1.00 -1.00 -1
13 | [2,] -10 -1.00 0.10 -3 -1.00 -1.00 -1
14 | [3,] -1  0.01 -1.00 -3  0.01 -1.00 -1
15 | [4,] -1 -1.00 0.01 -1  0.01  0.01 -1
16 | [5,] -1 -1.00 -1.00 -3 -1.00  0.01 100
17 | [6,] -1 -1.00 -1.00 -1  0.01 -1.00 100
18 | [7,] -1 -1.00 -1.00 -1 -1.00  0.01 100
```

Here's how you read this: the rows represent where you've come FROM, and the columns represent where you're going TO. Each element 1 through 7 corresponds directly to S1 through S7 in the cartoon above. Each cell contains a reward (or penalty, if the value is negative) if we arrive in that state.

The S1 state is *bad* for the robot... there's a hole in that piece of wood, so we'd really like to keep it away from that state. Location [1,1] on the matrix tells us what reward (or penalty) we'll receive if we start at S1 and *stay* at S1: -10 (that's bad). Similarly, location [2,1] on the matrix tells us that if we start at S2 and move left to S1, that's also bad and we should receive a penalty of -10. The S4 state is also undesirable –

there's a sticky patch there, so we'd like to keep the robot away from it. Location [3,4] on the matrix represents the action of going from S3 to S4 by moving right, which will put us on the sticky patch

Now load the qlearn command into your R session:

```

1 | qlearn <- function(R, N, alpha, gamma, tgt.state) {
2 |   # Adapted from https://stackoverflow.com/questions/39353580/how-to
3 |   Q <- matrix(rep(0,length(R)), nrow=nrow(R))
4 |   for (i in 1:N) {
5 |     cs <- sample(1:nrow(R), 1)
6 |     while (1) {
7 |       next.states <- which(R[cs,] > -1) # Get feasible actions
8 |       if (length(next.states)==1)       # There may only be one
9 |         ns <- next.states
10 |      else
11 |        ns <- sample(next.states,1) # Or you may have to pick from
12 |        if (ns > nrow(R)) { ns <- cs }
13 |        # NOW UPDATE THE Q-MATRIX
14 |        Q[cs,ns] <- Q[cs,ns] + alpha*(R[cs,ns] + gamma*max(Q[ns, ],
15 |        if (ns == tgt.state) break
16 |        cs <- ns
17 |      }
18 |    }
19 |    return(round(100*Q/max(Q)))
20 |  }

```

Run qlearn with the HOP rewards matrix, a learning rate of 0.1, a discount rate of 0.8, and a target state of S7 (the location to the far right of the wooden board). I did 10,000 episodes (where in each one, the robot dropped randomly onto the wooden board and has to get to S7):

```

1 | r.hop <- qlearn(HOP,10000,alpha=0.1,gamma=0.8,tgt.state=7)
2 | > r.hop
3 |   [,1] [,2] [,3] [,4] [,5] [,6] [,7]
4 | [1,]   0  51  64   0   0   0   0
5 | [2,]   0   0  64   0   0   0   0
6 | [3,]   0  51   0   0  80   0   0
7 | [4,]   0   0  64   0  80  80   0
8 | [5,]   0   0   0   0   0  80 100
9 | [6,]   0   0   0   0  80   0 100
10 | [7,]   0   0   0   0   0  80 100

```

The Q-Matrix that is presented encodes the best-value solutions from each state (the “policy”). Here’s how you read it:

- If you’re at s1 (first row), **hop** to s3 (biggest value in first row), then hop to s5 (go to row 3 and find biggest value), then hop to s7 (go to row 5 and find biggest value)
- If you’re at s2, go **right** to s3, then hop to s5, then hop to s7
- If you’re at s3, **hop** to s5, then hop to s7
- If you’re at s4, go **right** to s5 OR hop to s6, then go right to s7
- If you’re at s5, **hop** to s7
- If you’re at s6, go **right** to s7
- If you’re at s7, **stay there** (when you’re in the target state, the value function will not be able to pick out a “best action” because the best action is to do nothing)

Alternatively, the policy can be expressed as the best action from each of the 7 states: **HOP, RIGHT, HOP, RIGHT, HOP, RIGHT, (STAY PUT)**

Attempt 2: Use Reinforcement Learning Package

I also used the Reinforcement Learning package by Nicholas Proellocks (6/19/2017) described in <https://cran.r-project.org/web/packages/ReinforcementLearning/ReinforcementLearning.pdf> (<https://cran.r-project.org/web/packages/ReinforcementLearning/ReinforcementLearning.pdf>).

- **Input:** 1) a definition of the environment, 2) a list of states, 3) a list of actions, and 4) control parameters **alpha** (the learning rate; usually 0.1), **gamma** (the discount rate which describes how important future rewards are; often 0.9 indicating that 90% of the next reward will be taken into account), and **epsilon** (the probability that you'll try a random action; often 0.1)
- **Output:** A State-Action Value matrix, which attaches a number to how good it is to be in a particular state *and* take an action. You can use it to determine the highest value action from each state. (It contains the same information as the Q-matrix from Attempt 1, but you don't have to infer the action from the destination it brings you to.)
- **Pros:** Relatively straightforward. Allows you to specify epsilon, which controls the proportion of *random actions* you'll explore as you create episodes and explore your environment. **Cons:** Requires manual setup of all state transitions and associated rewards.

First, I created an "environment" that describes 1) how the states will change when actions are taken, and 2) what rewards will be accrued when that happens. I assigned a reward of -1 to all actions that are not special, e.g. landing on S1, landing on S4, or landing on S7. To be perfectly consistent with Attempt 1, I could have used 0.01 instead of -1, but the results will be similar. The values you choose for rewards are sort of arbitrary, but you do need to make sure there's a comparatively large positive reward at your target state and "negative rewards" for states you want to avoid or are physically impossible.

```

1 my.env <- function(state,action) {
2   next_state <- state
3   if (state == state("s1") && action == "right") { next_state <-
4   if (state == state("s1") && action == "hop") { next_state <-
5
6   if (state == state("s2") && action == "left") {
7     next_state <- state("s1"); reward <- -10 }
8   if (state == state("s2") && action == "right") { next_state <-
9   if (state == state("s2") && action == "hop") {
10    next_state <- state("s4"); reward <- -3 }
11
12   if (state == state("s3") && action == "left") { next_state <-
13   if (state == state("s3") && action == "right") {
14     next_state <- state("s4"); reward <- -3 }
15   if (state == state("s3") && action == "hop") { next_state <-
16
17   if (state == state("s4") && action == "left") { next_state <-
18   if (state == state("s4") && action == "right") { next_state <-
19   if (state == state("s4") && action == "hop") { next_state <-

```

```

20
21   if (state == state("s5") && action == "left") {
22     next_state <- state("s4"); reward <- -3 }
23   if (state == state("s5") && action == "right") { next_state <-
24   if (state == state("s5") && action == "hop") {
25     next_state <- state("s7"); reward <- 10 }
26
27   if (state == state("s6") && action == "left") { next_state <-
28   if (state == state("s6") && action == "right") {
29     next_state <- state("s7"); reward <- 10 }
30
31   if (next_state == state("s7") && state != state("s7")) {
32     reward <- 10
33   } else {
34     reward <- -1
35   }
36   out <- list(NextState = next_state, Reward = reward)
37   return(out)
38 }

```

Next, I installed and loaded up the ReinforcementLearning package and ran the RL simulation:

```

1  install.packages("ReinforcementLearning")
2  library(ReinforcementLearning)
3  states <- c("s1", "s2", "s3", "s4", "s5", "s6", "s7")
4  actions <- c("left", "right", "hop")
5  data <- sampleExperience(N=3000, env=my.env, states=states, actions=actions)
6  control <- list(alpha = 0.1, gamma = 0.8, epsilon = 0.1)
7  model <- ReinforcementLearning(data, s = "State", a = "Action", r = "Reward",
8    s_new = "NextState", control = control)

```

Now we can see the results:

```

1  > print(model)
2  State-Action function Q
3      hop      right      left
4  s1  2.456741  1.022440  1.035193
5  s2  2.441032  2.452331  1.054154
6  s3  4.233166  2.469494  1.048073
7  s4  4.179853  4.221801  2.422842
8  s5  6.397159  4.175642  2.456108
9  s6  4.217752  6.410110  4.223972
10 s7 -4.602003 -4.593739 -4.591626
11
12 Policy
13      s1      s2      s3      s4      s5      s6      s7
14  "hop" "right" "hop" "right" "hop" "right" "left"
15
16 Reward (last iteration)
17 [1] 223

```

The recommended policy is: **HOP, RIGHT, HOP, RIGHT, HOP, RIGHT, (STAY PUT)**

If you tried this example and it didn't produce the same response, don't worry! Model-free reinforcement learning is done by simulation, and when you used the sampleExperience function, you generated a different set of state transitions to learn from. You may need more samples, or to tweak your rewards structure, or both.)



5 replies »

🔗 Pingback: Reinforcement Learning: Q-Learning with the Hopping Robot – Mubashir Qasim

Hi,

Reply

Looks like there's a tiny typo in your R code (probalby copy-paste errors).

In attempt 2: "y.env" probalby is "my.env"

Great article.

Thank you! WordPress seems to chomp at least three parts of my code ' time I publish a post. I caught the first two, so thanks for catching the thiru 🙄

Reply

🔗 Pingback: Distilled News | Data Analytics & R

Thank you for your good example. I learn a lot more clearer about RL now.

Reply

Just in the Attempt 2 my.env function, it seems that the last "if" clause will overwrite the -3 penalty for state 4.

As a result, RL sometimes gets the results of good policy as "hop" for state 2 into state 4 without penalty faster towards state 7.

If changing s4 as -3 for the data, RL can get the correct results.

Thank you.