

Optimizing adaptive profiling in Alleria using Reinforcement Learning

Shreshth Tuli

Department of CSE, IIT Delhi

Abstract

Alleria is a new framework for generating instruction and detailed memory traces. It can be used by researchers to collect interesting information about one or more target applications. The earlier proposed heuristic based adaptive profiling mechanism has great performance improvement with respect to its counterparts. This document throws light on one improvement that is made to Alleria using Reinforcement Learning algorithms. It makes Alleria intelligent in the sense that it can figure out by itself how many processing threads and buffers are required and dynamically adjust these parameters. This intelligence is provided using Q -Learning algorithm with ϵ -greedy approach.

1 Introduction : Q -Learning

Q -learning is a reinforcement learning algorithm that tries to find optimal actions by learning a state-action value function. The state-action value function, or simply $Q(s, a)$, is a table having rows as states, actions as columns, and values as entries. Values quantify the reward that the agent is expected to collect if it selects the corresponding action from the corresponding state and proceed from-there-on-out by following the optimal policy. Thus if the value function is known, then the optimal policy is simply to select the action having the highest value for the current state. The issue is that the Q -values are initially unknown, so they have to be learned in some way. Q -learning learns values by exploring the state-action space (literally just trying different action) and updating Q based on rewards that are measured after each action.

Thus, the two important components of Q -learning are:

1. Exploring state-action space.
2. Updating Q

Many flavors of algorithms have been devised to address these points; we'll focus on one of the simpler and more widely used ones.

1.1 Exploring state-action space

Exploring state-action space requires us to visit several state-actions and measure their instantaneous rewards, i.e., the reward obtained by performing action a from state s . Generally state-action space is huge, even for simple problems, so it's necessary to focus exploration in promising regions. We use an ϵ -greedy approach where the optimal action ($\arg \max_a Q_t(s, a)$) is selected with probability $1 - \epsilon$, while an arbitrary action is selected with probability ϵ . In addition, we apply a decay that decreases ϵ over time. This causes the agent to exploit information it's already collected, instead of continuing to explore un-visited state-actions.

1.2 Updating Q

Q is updated according to the following equation:

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha(r + \gamma(\max_a Q_t(s', a) - Q_t(s, a)))$$

s is the agent's current state, a is the action it performs, r is the instantaneous reward, s' is the state the agent finds itself in after the action is performed, and γ and α are parameters.

The update equation is identical to an on-the-fly average $x_{t+1}^- = \bar{x}_t + 1/t(x_t - \bar{x}_t)$ where the sample is $r + \gamma \max_a Q_t(s', a)$. This sample combines the current reward r with the future reward $\max_a Q_t(s', a)$ to let information leak backward from s' to s . This combination allows strong instantaneous rewards to be lessened if the future state is bad, which is intuitively appropriate. The factor γ determines the influence of future values on current values, and α plays the role of a learning rate.

2 Implementation

The implementation of the model with ϵ -greedy algorithm has been done on C++.

The states include the possible pairs (θ_1, θ_2) , where θ_1 denotes the number of thread counts and θ_2 denotes the number of buffers. The performance parameter denoted by P is used directly as reward.

The procedural approach can be translated into plain English steps as follows:

1. Initialize the Q-values table, $Q(s, a)$.
2. Observe the current state, s .
3. Choose an action, a , for that state based on one of the action selection policies explained here on the previous page
4. Take the action, and observe the reward, r , as well as the new state, s' .
5. Update the Q-value for the state using the observed reward and the maximum reward possible for the next state. The updating is done according to the formula and parameters described above.
6. Set the state to the new state, and repeat the process until a terminal state is reached.

Appendices

C++ Code

```
1  /*
2  Author: Shreshth Tuli
3  Date: 26/02/2018
4
5  This program functions as the "brain" of a reinforcement
6  learning
7  agent whose goal is to provide the optimal values of number of
8  buffers and
9  number of threads for efficient profiling.
10 For use in Alleria profiler code
11 */
12
13 #include <iostream>
14 #include <random.h>
15
16 ///////////////////////////////////////////////////
17 //Computational parameters
18 float gamma = 0.75; //look-ahead weight
19 float alpha = 0.2; // "Forgetfulness" weight. The closer
20 // this is to 1 the more weight is given to recent samples.
21 // A high value is kept because of a
22 // highly dynamic situation, we cannot keep it very high as then
23 // the system might not converge
24
25 //Parameters for getAction()
26 float epsilon; //epsilon is the probability of
27 // choosing an action randomly. 1-epsilon is the probability of
28 // choosing the optimal action
29
30 //Variable 1 Parameters - Number of threads
31 const int numTheta1States = 10;
32 float theta1InitialCount = 1;
33 float theta1Max = 10;
34 float theta1Min = 1;
35 float deltaTheta1 = 1;
36 int s1 = int((theta1InitialCount - theta1Min)/deltaTheta1);
37 //This is an integer between zero and
38 // numTheta1States-1 used to index the state number of servo1
```

```

32 //Variable 2 Parameters – Number of Buffers
33 const int numTheta2States = 10;
34 float theta2InitialCount = 1;
35 float theta2Max = 10;
36 float theta2Min = 1;
37 float deltaTheta2 = 1;
38 int s2 = int((theta2InitialCount - theta2Min)/deltaTheta2);
           //This is an integer between zero and
           numTheta2States-1 used to index the state number of servo2
39
40 //Initialize Q to zeros
41 const int numStates = numTheta1States*numTheta2States;
42 const int numActions = 5;
43 float Q[numStates][numActions];
44
45 //Initialize the state number. The state number is calculated
           using the theta1 state number and
46 //the theta2 state number. This is the row index of the state
           in the matrix Q. Starts indexing at 0.
47 int s = int(s1*numTheta2States + s2);
48 int sPrime = s;
49
50 //Initialize vars for getDeltaDistance()
51 float distanceNew = 0.0;
52 float distanceOld = 0.0;
53 float deltaDistance = 0.0;
54
55 //These get used in the main loop
56 float r = 0.0;
57 float lookAheadValue = 0.0;
58 float sample = 0.0;
59 int a = 0;
60
61 //Returns an action 0, 1, ... 4 : NONE, theta1++, theta1--,
           theta2++, theta2--
62 int getAction(){
63     int action;
64     float valMax = -10000000.0;
65     float val;
66     int aMax;
67     float randVal;
68     int allowedActions[5] = {-1, -1, -1, -1, -1}; // -1 if action
           of the index takes you outside the state space. +1 otherwise
69     boolean randomActionFound = false;
70

```

```

71 //find the optimal action. Exclude actions that take you
    outside the allowed-state space.
72 if((s1 + 1) != numTheta1States){
73     allowedActions[0] = 1;
74     val = Q[s][0];
75     if(val > valMax){
76         valMax = val;
77         aMax = 0;
78     }
79 }
80 if(s1 != 0){
81     allowedActions[1] = 1;
82     val = Q[s][1];
83     if(val > valMax){
84         valMax = val;
85         aMax = 1;
86     }
87 }
88 if((s2 + 1) != numTheta2States){
89     allowedActions[2] = 1;
90     val = Q[s][2];
91     if(val > valMax){
92         valMax = val;
93         aMax = 2;
94     }
95 }
96 if(s2 != 0){
97     allowedActions[3] = 1;
98     val = Q[s][3];
99     if(val > valMax){
100         valMax = val;
101         aMax = 3;
102     }
103 }
104
105 //implement epsilon greedy
106 randVal = float(random(0,101));
107 if(randVal < (1.0-epsilon)*100.0){ //choose the optimal
    action with probability 1-epsilon
108     action = aMax;
109 }else{
110     while(!randomActionFound){
111         action = int(random(0,5)); //otherwise pick an
            action between 0 and 4 randomly (inclusive), but don't use
            actions that take you outside the state-space

```

```

112         if(allowedActions[action] == 1){
113             randomActionFound = true;
114         }
115     }
116 }
117 return(action);
118 }
119
120 //Given a and the global(s) find the next state. Also keep
121 //track of the individual joint indexes s1 and s2.
122 void setSPrime(int action){
123     if(action == 0){
124         //NONE
125         sPrime = s
126     }
127     else if (action == 1){
128         //theta1++
129         sPrime = s + numTheta2States;
130         s1++;
131     } else if (action == 2){
132         //theta1--
133         sPrime = s - numTheta2States;
134         s1--;
135     } else if (action == 3){
136         //theta2++
137         sPrime = s + 1;
138         s2++;
139     } else {
140         //theta2--
141         sPrime = s - 1;
142         s2--;
143     }
144 }
145
146 //Update the number of threads and buffers (this is the physical
147 //state transition command)
148 void setPhysicalState(int action){
149     float currentCount;
150     float finalCount;
151     if (action == 1){
152         currentCount = //theta 1 read
153         finalCount = currentCount + deltaTheta1;
154         //theta1 write finalCount
155     } else if (action == 2){

```

```

155     currentCount = //theta 1 read
156     finalCount = currentCount - deltaTheta1;
157     //theta1 write finalCount
158 }else if (action == 3){
159     currentCount = //theta2 read
160     finalCount = currentCount + deltaTheta2;
161     //theta2 write finalCount
162 }else if (action == 4){
163     currentCount = //theta2 read
164     finalCount = currentCount - deltaTheta2;
165     //theta2 write finalCount
166 }
167 }
168
169
170 //Get the reward using the increase in performance since the
    last call
171 float getDeltaDistance() {
172     //get current performance
173     distanceNew = // read current performance
174     deltaDistance = distanceNew - distanceOld;
175     //if (abs(deltaDistance) < 57.0 || abs(deltaDistance) > 230.0)
        {
176         //don't count noise
177         // deltaDistance = 0.0;
178     }
179     distanceOld = distanceNew;
180     return deltaDistance;
181 }
182
183 //Get max over a' of Q(s',a'), but be careful not to look at
    actions which take the agent outside of the allowed state
    space
184 float getLookAhead() {
185     float valMax = -100000.0;
186     float val;
187     if ((s1 + 1) != numTheta1States) {
188         val = Q[sPrime][0];
189         if (val > valMax) {
190             valMax = val;
191         }
192     }
193     if (s1 != 0) {
194         val = Q[sPrime][1];
195         if (val > valMax) {

```



```

196     valMax = val;
197 }
198 }
199 if((s2 + 1) != numTheta2States){
200     val = Q[sPrime][2];
201     if(val > valMax){
202         valMax = val;
203     }
204 }
205 if(s2 != 0){
206     val = Q[sPrime][3];
207     if(val > valMax){
208         valMax = val;
209     }
210 }
211 return valMax;
212 }
213
214 void initializeQ(){
215     for(int i=0; i<numStates; i++){
216         for(int j=0; j<numActions; j++){
217             Q[i][j] = 10.0; //Initialize to a positive
                             //number to represent optimism over all state-actions
218         }
219     }
220 }
221
222 ///////////////////////////////////////////////////
223
224 const int readDelay = 200; //allow time for
                             //the agent to execute after it sets its physical state
225 const float explorationMinutes = 1.0; //the desired
                             //exploration time in minutes
226 const float explorationConst = (explorationMinutes*60.0)/((float
                             //this is the approximate exploration
                             //time in units of number of times through the loop
                             (readDelay))/1000.0);
227
228 int t = 0;
229 void main(){
230     while (true){
231         t++;
232         epsilon = exp(-float(t)/explorationConst);
233         a = getAction(); //a is between 0 and 4
234         setSPrime(a); //this also updates s1 and s2.
235         setPhysicalState(a);

```

```

236     delay(readDelay); //put a delay after
the physical action occurs so the agent has some delay
between two performance reads
237     r = getDeltaDistance();
238     lookAheadValue = getLookAhead();
239     sample = r + gamma*lookAheadValue;
240     Q[s][a] = Q[s][a] + alpha*(sample - Q[s][a]);
241     s = sPrime;
242
243     if(t == 2){ //need to reset Q at the
beginning since a spurious value may arise at the first
initialization
244         initializeQ();
245     }
246 }
247 }

```

Listing 1: C++ code

References

- [1] <http://andraugust.com/rl/>
- [2] <https://www.cse.unsw.edu.au/~cs9417ml/RL1/algorithms.html>
- [3] <https://towardsdatascience.com/introduction-to-various-reinforcement-learning-algorithms-i-q-learning-sarsa-dqn-ddpg-72a5e0cb6287>
- [4] <https://sites.ualberta.ca/~szepesva/papers/RLAlgsInMDPs.pdf>
- [5] https://en.wikipedia.org/wiki/Reinforcement_learning
- [6] <https://github.com/samindaa/RLLib>