

# Stochastic Optimization Assignment 3

Aytek Mutlu

December 09, 2019

## 1 Introduction

In this assignment, 2 Bernoulli arms with uniform priors are considered. The objective is to find the maximum discounted number of successes using those arms. For this purpose, variety of different policies are applied and compared in the following sections.

## 2 Gittins Index

In order to apply Bayesian policy, Gittins Indices are calculated for all possible states offline. Considering the expensive computation of the mentioned index, the state space is limited to (40,40), which means the maximum number of success and failure is bound to 40. As required in the assignment, discount factor is used as 0.95 i.e.  $\alpha = 0.95$

For each and every state in state space, an iteration is looped until value matrix converges (i.e. change in consecutive iterations are less than  $10^{-3}$ ). Within this iteration, value of each and every state is updated with the following formula:

$$V_t[k, l] = \max(p + \alpha \cdot (p \cdot V_{t-1}[k+1, l] + (1-p) \cdot V_{t-1}[k, l+1]), V_{t-1}[a, b])$$

where  $p = k/(k+l)$  and state  $(a, b)$  is the state where restarting option applies to. This means  $(a, b)$  is the state for which Gittins Index is calculated.

When  $V_t[a, b] - V_{t-1}[a, b] < 0.001$ , iteration stops and Gittins index for state  $(a, b)$  is calculated as follows:

$$m(a, b) = V[a, b] \cdot (1 - \alpha)$$

This procedure is applied for all states in state space, and the following Gittins index table is obtained (Only first 10 successes are 10 failures are illustrated for visual concerns):

	k	1	2	3	4	5	6	7	8	9	10
l											
1		0.7583	0.8354	0.8710	0.8924	0.9069	0.9175	0.9256	0.9321	0.9374	0.9418
2		0.5570	0.6776	0.7409	0.7812	0.8095	0.8308	0.8473	0.8605	0.8715	0.8807
3		0.4313	0.5591	0.6359	0.6868	0.7246	0.7533	0.7762	0.7947	0.8102	0.8234
4		0.3465	0.4731	0.5527	0.6101	0.6530	0.6864	0.7138	0.7362	0.7550	0.7714
5		0.2872	0.4078	0.4875	0.5465	0.5927	0.6294	0.6593	0.6845	0.7061	0.7246
6		0.2440	0.3567	0.4355	0.4941	0.5413	0.5800	0.6120	0.6391	0.6622	0.6824
7		0.2112	0.3169	0.3926	0.4510	0.4975	0.5369	0.5703	0.5987	0.6233	0.6447
8		0.1857	0.2843	0.3567	0.4141	0.4609	0.4994	0.5333	0.5627	0.5882	0.6106
9		0.1653	0.2576	0.3269	0.3826	0.4287	0.4675	0.5005	0.5303	0.5565	0.5796
10		0.1488	0.2353	0.3015	0.3552	0.4005	0.4390	0.4723	0.5012	0.5278	0.5514

Table 1: Gittins Indices for Bernoulli Arms with k success and l failures

### 3 Policies

For the sake of clear comparison, all policies are implemented with same sample space of  $N = 40$ . Each policy calculates the discounted number of successes and in order to come up with the expected discounted number of successes, all policies are run for  $N_{simulation} = 1000$  times and results are averaged. Then, the results are compared in Section 4.

#### 3.1 Full Information

This policy is the base policy where it is assumed that the priors of bandits are known and therefore better arm is selected all the time. Of course, a different arm with a different given prior is selected at each simulation. Here are the steps for the policy:

1. Pick two priors for 2 arms (either given manually or use random numbers between 0 and 1)
2. Select the arm with larger prior and use the same selected arm for 40 times
3. Simulate each time based on given prior of the selected arm (i.e.  $x = 1$  if trial is a success and vice versa)
4. Update reward with  $r = r + \alpha^t \cdot x$

#### 3.2 Bayesian using Gittins Index

Bayesian policy bases its methodology on Gittins Index table calculated offline. Conceptually, it selects the arm with higher Gittins index given its state at each iteration. Here are the steps for the policy:

1. Randomly pick an initial arm
2. Simulate selected arm (initially with 0.5 success probability, and then use posterior distribution in further iterations i.e. no. of successes / no. of trials)

3. Update state of the selected arm (i.e. 1 success 0 failure etc.)
4. Update reward with  $r = r + \alpha^t \cdot x$
5. Compare Gittins indices of both arms given their states and then select the arm with higher Gittins index
6. Go to step 2 and iterate for 40 times

### 3.3 Thompson Sampling

Thompson sampling follows a very similar procedure with Bayesian approach where the only difference is that the arm at each iteration is selected based on their posterior Beta distributions. Here are the steps:

1. Randomly pick an initial arm
2. Simulate selected arm (initially with 0.5 success probability, and then use posterior distribution in further iterations i.e. no. of successes / no. of trials)
3. Update state of the selected arm (i.e. 1 success 0 failure etc.)
4. Update reward with  $r = r + \alpha^t \cdot x$
5. Compare posterior Beta distributions of both arms given their states and then select the arm with higher Beta distribution
6. Go to step 2 and iterate for 40 times

### 3.4 Greedy Q-Learning

Greedy Q-Learning policy is a stateless approach where the policy always exploits the arm with larger Q value. Here are the steps for Greedy Q-Learning:

1. Initialize Q values with 0 for both arms
2. Provide 2 priors for arms
3. Select the arm with larger Q-value
4. Simulate selected arm with its prior
5. Update reward with  $r = r + \alpha^t \cdot x$
6. Update Q value of the selected arm with  $Q[arm] = (1/t) \cdot \alpha^t \cdot x + (1 - (1/t)) \cdot Q[arm]$
7. Go to step 3 and iterate for 40 times

### 3.5 Optimistic Q-Learning

Optimistic Q-learning is exactly same with Greedy Q-Learning where the only difference is the initial values for Q. Higher values are used during initialization and the rest of the process is identical.

1. Initialize Q values with 1 for both arms
2. Provide 2 priors for arms
3. Select the arm with larger Q-value
4. Simulate selected arm with its prior
5. Update reward with  $r = r + \alpha^t \cdot x$
6. Update Q value of the selected arm with  $Q[arm] = (1/t) \cdot \alpha^t \cdot x + (1 - (1/t)) \cdot Q[arm]$
7. Go to step 3 and iterate for 40 times

### 3.6 $\epsilon$ -Greedy Q-Learning

$\epsilon$ -Greedy Q-Learning explores with  $\epsilon$  probability and exploits with the rest of the probability. This means that policy picks a random arm during exploration and picks the arm that has maximum Q value during exploitation.

1. Initialize Q values with 0 for both arms
2. Provide 2 priors for arms
3. Select an arm either randomly with  $\epsilon$  probability or the arm with larger Q-value with  $1 - \epsilon$  probability
4. Simulate selected arm with its prior
5. Update reward with  $r = r + \alpha^t \cdot x$
6. Update Q value of the selected arm with  $Q[arm] = (1/t) \cdot \alpha^t \cdot x + (1 - (1/t)) \cdot Q[arm]$
7. Go to step 3 and iterate for 40 times

## 4 Conclusion

Only parameter that would affect the outcome of any policy is  $\epsilon$  in  $\epsilon$ -Greedy Q-Learning. For this purpose, an experiment is run with different  $\epsilon$  values and Figure 1 is obtained.

As it can be seen,  $\epsilon = 0.3$  is the optimal point for the objective. The comparison of policies are shown in Table 2. Not surprisingly, Full Information provided the maximum number, as it treats on the provided information.

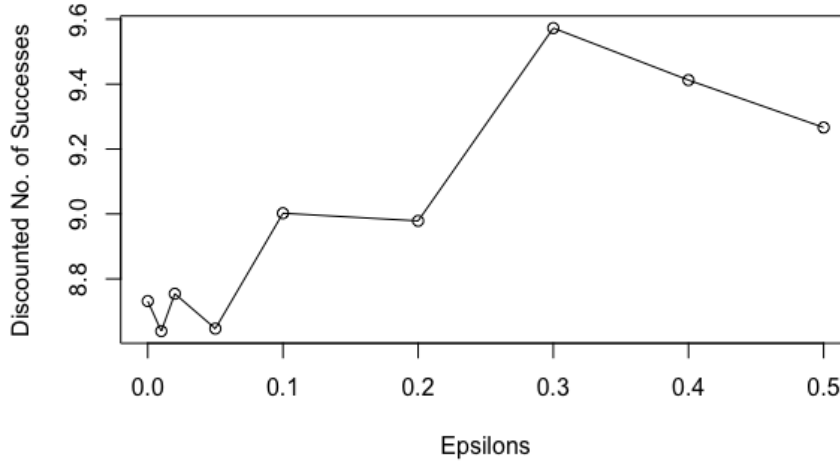


Figure 1: Discounted Number of Successes with changing  $\epsilon$  values

Bayesian policy with Gittins Index follows the full information, yet its computational hassle makes it less favorable. Then, Thompson Sampling is very successful and computationally not heavy, as well.

Among stateless approaches, Greedy Q-learning is the worst as it focuses too much on exploitation. Its optimistic version offers a better result, as it is initialized with higher values. Finally,  $\epsilon$ -Greedy Q-Learning performs moderately well as it looks for a balance between exploration and exploitation.

Policy	Discounted No. of Successes
Full Information	11.58
Bayesian with Gittins Index	10.72
Thompson Sampling	10.08
Greedy Q-Learning	8.51
Optimistic Q-Learning	9.79
$\epsilon$ -Greedy Q-Learning	9.28

Table 2: Objective Results of Policies

## 5 Code

```
1
2
3 #Gittins index calculation function
4 Gittins ← function(alpha,N){
5   m = matrix(c(0),N,N)
6   for(a in 1:N){
7     for(b in 1:N){
8       V = matrix(c(0),N,N)
9       diff = 1.0
10      while(diff > 0.001){
11        V_temp = V
12        for(k in a:N){
13          for(l in b:N){
14            p = (k)/(k+1)
15            V[k,l] = max(p + alpha * (p * V_temp[min(k+1,N),
16              l] + (1-p) * V_temp[k,min(l+1,N)]), V_temp[a,
17                b])
18          }
19        }
20        diff = V[a,b] - V_temp[a,b]
21      }
22      m[a,b] = V[a,b]*(1-alpha)
23    }
24  }
25
26 FullInformation ← function(alpha,N){
27   p_bandits = runif(2)
28   arm = which.max(p_bandits)
29   r=0
30   for(t in 1:N){
31     p_arm = p_bandits[arm]
32     x = rbinom(1,1,p_arm)
33     r = r + alpha^(t-1) *x
34   }
35   return(r)
36 }
37
38 #Bayesian with Gittins Index
39 Bayesian ← function(alpha,N){
40   arm_states = matrix(c(0),2,2)
41   p_bandits = runif(2)
42   arm = which.max(p_bandits)
43   r=0
44   for(t in 1:N){
45     p_arm = (arm_states[arm,2]+1)/(rowSums(arm_states)[arm
```

```

    ]+2)
46     x = rbinom(1,1,p_arm)
47     arm_states[arm,x+1] = arm_states[arm,x+1]+1
48     r = r + alpha^(t-1) *x
49     arm = ifelse(gittins_table[arm_states[1,2]+1,arm_states
      [1,1]+1] > gittins_table[arm_states[2,2]+1,arm_states
      [2,1]+1],1,2)
50   }
51   return(r)
52 }
53
54 #Thompson Sampling
55 Thompson <- function(alpha,N){
56   arm_states = matrix(c(0),2,2)
57   p_bandits = runif(2)
58   arm = which.max(p_bandits)
59   r=0
60   for(t in 1:N){
61     p_arm = (arm_states[arm,2]+1)/(rowSums(arm_states)[arm
      ]+2)
62     x = rbinom(1,1,p_arm)
63     arm_states[arm,x+1] = arm_states[arm,x+1]+1
64     r = r + alpha^(t-1) *x
65     arm = ifelse(rbeta(n=1,shape1=arm_states[1,2]+1,shape2 =
      arm_states[1,1]+1)>rbeta(n=1,shape1=arm_states
      [2,2]+1,shape2 =arm_states[2,1]+1),1,2)
66   }
67   return(r)
68 }
69
70 #Greedy Q-Learning
71 Greedy <- function(alpha,N,epsilon,initial_prob){
72
73   Q = matrix(c(initial_prob),1,2)
74   p_bandits = runif(2)
75   r=0
76
77   for(t in 1:N){
78
79     if(runif(1)<epsilon)
80       arm = rbinom(1,1,0.5)+1
81     else
82       arm = which.max(Q)
83
84     x = rbinom(1,1,p_bandits[arm])
85     r = r + alpha^(t-1) *x
86
87     Q[arm] = (1/t) * alpha^(t-1) * x + (1-(1/t))*Q[arm]
88   }
89   return(r)

```

```

90 }
91
92 #calculate Gittins Index
93 alpha = 0.95
94 N = 40
95 gittins_table = Gittins(alpha,N)
96
97 #Simulation
98 N_sim = 1000
99 full_info_reward = 0
100 bayesian_reward = 0
101 thompson_reward = 0
102 greedy_reward = 0
103 eps_greedy_reward = 0
104 opt_greedy_reward = 0
105
106 for(n in 1:N_sim){
107
108     fi = FullInformation(alpha,N-1)
109     ba = Bayesian(alpha,N-1)
110     th = Thompson(alpha,N-1)
111     gr = Greedy(alpha,N-1,epsilon=0,initial_prob = 0)
112     opt_gr = Greedy(alpha,N-1,epsilon=0,initial_prob = 1)
113     eps_gr = Greedy(alpha,N-1,epsilon=0.3,initial_prob = 0)
114
115
116     full_info_reward = full_info_reward + fi
117     bayesian_reward = bayesian_reward + ba
118     thompson_reward = thompson_reward + th
119     greedy_reward = greedy_reward + gr
120     opt_greedy_reward = opt_greedy_reward + opt_gr
121     eps_greedy_reward = eps_greedy_reward + eps_gr
122 }
123
124 full_info_reward = full_info_reward/N_sim
125 bayesian_reward = bayesian_reward/N_sim
126 thompson_reward = thompson_reward/N_sim
127 greedy_reward = greedy_reward/N_sim
128 opt_greedy_reward = opt_greedy_reward/N_sim
129 eps_greedy_reward = eps_greedy_reward/N_sim
130
131 print(c("Full info:",full_info_reward))
132 print(c("Bayesian:",bayesian_reward))
133 print(c("Thompson Sampling:",thompson_reward))
134 print(c("Greedy Q-Learning:",greedy_reward))
135 print(c("Optimistic Greedy Q-Learning:",opt_greedy_reward))
136 print(c("Epsilon-Greedy Q-Learning:",eps_greedy_reward))
137
138 #Find optimal epsilon
139 epsilons = c(0,0.01,0.02,0.05,0.1,0.2,0.3,0.4,0.5)

```



```

140 eps_greedy_rewards = c()
141
142 for(eps in epsilons){
143     eps_greedy_reward = 0
144     for(n in 1:N_sim){
145         eps_gr = Greedy(alpha,N-1,epsilon=eps,initial_prob = 0)
146         eps_greedy_reward = eps_greedy_reward + eps_gr
147     }
148     eps_greedy_rewards = c(eps_greedy_rewards,
149                             eps_greedy_reward)
150 }
151 eps_greedy_rewards = eps_greedy_rewards/N_sim
152
153 plot(epsilons,eps_greedy_rewards,type='o',xlab='Epsilons',
154      ylab='Discounted No. of Successes')

```