

Multi-armed Bandits with externalities

VVN

2019

Reward structure :

- We have a set of N arms and a corresponding set of N user types.
- If user is of type i and is shown arm j , then the reward structure is as follows :
 - If $i \neq j$ then $\mathbf{E}(R_t) = \mu_{ij}$ with reward $R_t \in \{-1, +1\}$.
 - If $i = j$ then $\mathbf{E}(R_t) = \mu_{ii}$ with reward $R_t \in \{0, 2\}$.
- Therefore we have a matrix $M = [\mu_{ij}]_{N \times N}$ of the mean rewards where in each row i , the element μ_{ii} has the highest value.

Updating user preferences :

- We have a user preference vector α with N components with:

$$Prob(\text{User arrives with type } i) = \frac{\alpha_i}{\sum_j \alpha_j}$$

- If A_p is the user preference and A_r is the recommended arm with R_t being the reward accrued, then :
 - $\alpha_{A_p} - = R_t$
 - $\alpha_{A_r} + = R_t$
- Thus the rewards and the policy together control the structure of the population α .
- The arm preference of the user is revealed to the recommender after he has chosen the arm to be shown and has obtained the reward.

The Aim

- We define the "best arm" to be the arm $\operatorname{argmax}_i(\mu_{ii})$
- We aim to manipulate the population vector α through our policy so that the probability of a user arriving with preference for the best arm is maximized.
- If it is not possible to bring about the desired α , we also provide the minimal conditions on the mean matrix M so that our policy works.
- A secondary aim might be to minimize the cumulative regret accrued in the process of achieving the desired α .

Our first test policy is to choose all arms uniformly at random. After simulating this policy many times and observing the trend for average α for fixed matrix M , we make the following observations :

- The values of α_i 's arrange themselves in the order of the column sums of the matrix M .
- Even if the arm i is the best arm, α_i can drop to 0 very quickly if the i^{th} column sum is lower than that for other arms.
- The α trajectory for such a case is shown in the following figure.

Test policy 1 simulated

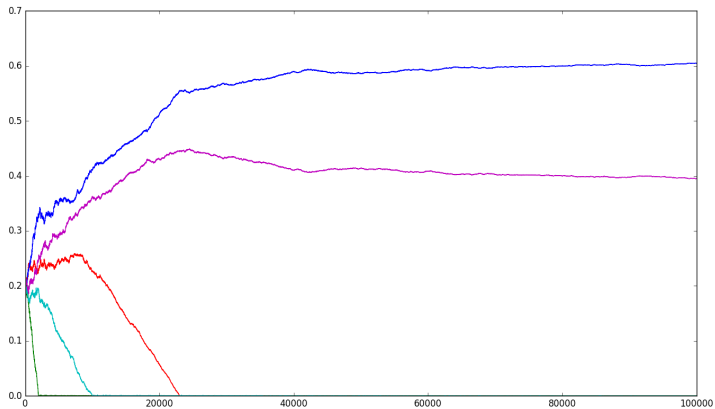
- The matrix we used is the following:

$$\begin{bmatrix} 0.7 & 0.1 & 0.3 & 0.3 & 0.4 \\ 0.1 & 0.9 & 0.4 & 0.4 & 0.4 \\ 0.2 & 0.05 & 0.6 & 0.4 & 0.5 \\ 0.6 & 0.2 & 0.2 & 0.65 & 0.2 \\ 0.2 & 0.2 & 0.55 & 0.2 & 0.55 \end{bmatrix}$$

- The color scheme is in the order : (r,g,b,c,m).

Test policy 1 simulated

Figure 1: Uniform policy



Our second test policy is to choose the same one arm repeatedly. The following were the observations after the simulation:

- If we are repeatedly choosing arm i then α_j keeps reducing if μ_{ji} is positive.
- If an arm j has a large column sum (relative to other arms) has no positive elements in its row, then it is not possible to reduce α_j using any particular arm.

Based on the observations obtained after implementing the previous two policies, we now implement the following algorithm:

- **Step 1 - Estimation phase :**

- We first need to estimate the matrix M so as to decide the best arm and then start manipulating the α_i 's.
- For the first T iterations, we select arms uniformly at random and update our estimate \hat{M} for the matrix.
- The key tradeoff involved here is that if we take too long to estimate the matrices, then the α_i of the best arm may fall to zero by virtue of a low column sum.
- Therefore choosing an appropriate value of T is essential, otherwise we would not be able to guess the best arm with surety, or its α_i would just drop down to zero while estimation.

- **Step 2 - Manipulation phase :**

- We now focus on an arm i such that:

$$i = \operatorname{argmax}_{j \neq \text{bestarm}}(\alpha_j)$$

Call this the target arm. This arm may change each with each iteration.

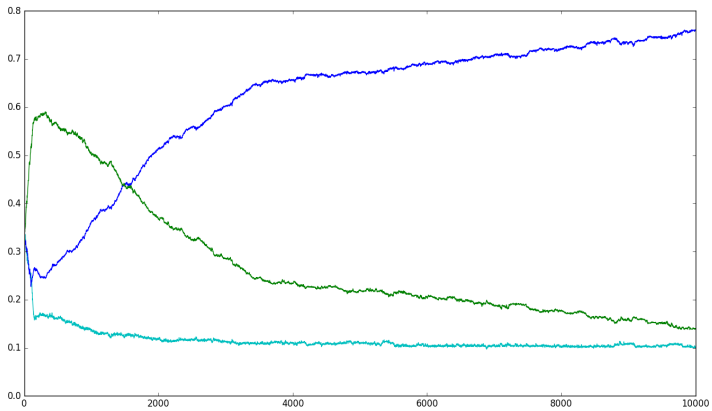
- Then we look at the i^{th} row in the estimated matrix \hat{M} and choose the arm k such that:

$$k = \operatorname{argmax}_{j \neq \text{bestarm}, j \neq i}(\mu_{ij})$$

- **Step 3 : Exploitation phase :**

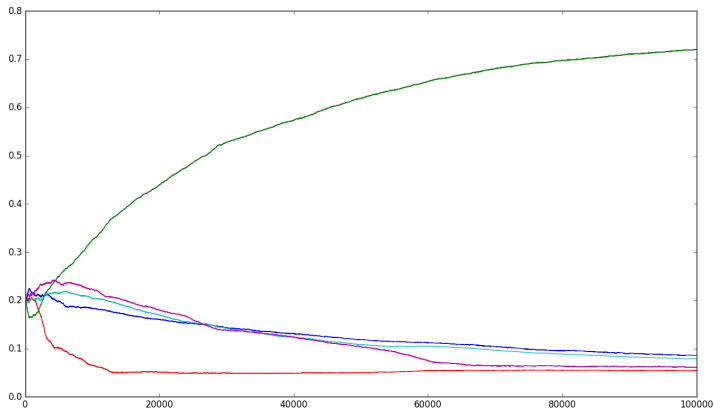
- If $\alpha_{\text{bestarm}} > 1 - \delta$ then we choose the best arm.

Figure 2: Policy with three arms



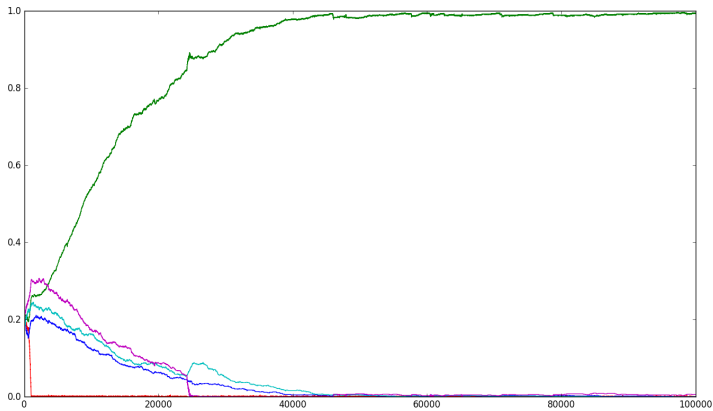
A Working policy : Simulations

Figure 3: Policy with five arms (averaged over 100 simulations)



A Working policy : Simulations

Figure 4: Policy with five arms (one particular simulation)

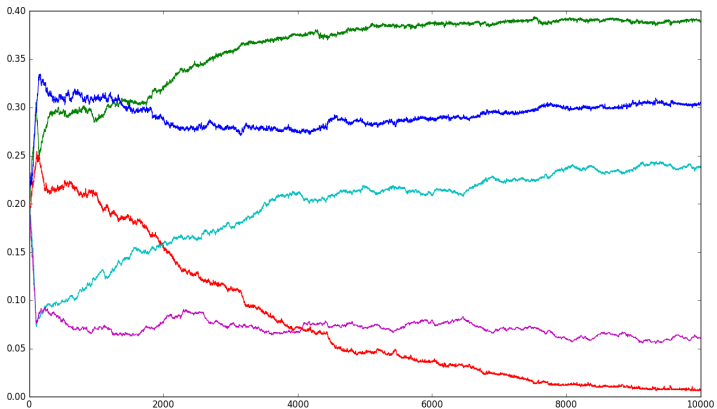


When does this not work?

- The policy that we discussed above does not work when there exists a row in the matrix M such that there is no positive mean reward in that row except for the diagonal element.
- In such cases there is no way to bring down the α value for this arm.
- Such a case is illustrated in the following simulation.

When does this not work?

Figure 5: Policy with five arms (not working)



Modifications to be done

- During the phase of the algorithm when the matrix M is being estimated, we need an algorithm that would do the estimation and at the same time keep the α values above zero with some tolerance.
- We also need to quantify the steady state α_i value that would be reached in presence of a "whiner" population (i.e the presence of a row in M with negative values of μ).
- If the column corresponding to a suboptimal arm is large then we need a way to ensure that the arm does not make its corresponding α go to 1 prematurely during the matrix estimation.

- *Here* is the github link for the simulation code.