

Multi-armed bandit problem

Anežka Lhotáková

June 25, 2022

Abstract

In probability theory and machine learning, multi-armed bandit problem refers to a model which can be seen as a set of real distributions $\{F_1, \dots, F_K\}$, each distribution being associated with the rewards delivered by one of the $K \in \mathbb{N}^+$ levers. Let μ_1, \dots, μ_K be the mean values associated with these reward distributions. The gambler iteratively plays one lever per round and observes the associated reward. The objective is to maximize the sum of the collected rewards. Several strategies or algorithms have been proposed as a solution to this problem in the last two decades. In this project, we will focus on the most popular strategies and their parameter optimization through heuristic approach.

Introduction

In this paper, model with i.i.d. (independent and identically distributed) rewards is presented. We define

- horizon H as the number of rounds to be played,
- maximum mean reward as $\mu^* = \max\{\mu_1, \mu_2, \dots, \mu_K\}$,
- regret ρ after T rounds played as

$$\rho(T) = T\mu^* - \sum_{t=1}^T r_t,$$

where r_t is the reward in round t (ρ is also a random variable),

- expected regret as $\mathbb{E}[\rho(T)]$,
- symbol ℓ_t^i , ($i = 1, \dots, K$ and $t = 1, \dots, H$) for the levers.

The protocol to the model is following:

Given H rounds to play, considering K levers, in each round $t = 1, \dots, H$:

1. Based on the strategy, the algorithm picks a lever ℓ_t^i , $i = 1, \dots, K$.
2. Algorithm observes reward r_t for the chosen lever.
3. If $t < H$, the algorithm will play again (point 1., $t \rightarrow t + 1$), otherwise the game is over.

The goal is to maximize the total reward over the T rounds. Above that, we make two important assumptions:

- The algorithm observes only the reward for the selected action, and nothing else.
- The reward for each action is i.i.d. For each lever ℓ^i , there is a distribution F_i over reals, called the reward distribution. Every time the lever is chosen, the reward is sampled independently from its distribution. The reward distributions are initially unknown to the algorithm.

For better understanding of the model, let us begin with an example. Imagine player walking into a casino with two slot-machines. Both slot-machines ℓ_1, ℓ_2 have their already given distributions F_1, F_2 . Let the reward r of each slot-machine be given as $r_1 \sim F_1 = \mathcal{N}(10, 3)$ and $r_2 \sim F_2 = \mathcal{N}(8, 4)$. Unfortunately, this information is hidden from the player. What possible strategies can the player use in order to maximize his win over $H = 100$ rounds? The mean optimal win is

$$r_O = H \cdot \mu^* = H \cdot \max\{\mu_1, \mu_2\} = 100 \cdot 10 = 1\,000.$$

a) *Explore only*: Firstly, he could choose to pull the levers randomly for the whole 100 rounds. In that case, his mean reward after 100 games would be

$$r = 50 \cdot 10 + 50 \cdot 8 = 900,$$

meaning expected regret is equal to $\mathbb{E}[\rho] = 1\,000 - 900 = 100$.

b) *Exploit only*: Second strategy could be based on first experience. One round on each slot-machine could be played and explored. In first round, player pulls lever ℓ_A and gets a reward of 9. Then, pulling the second lever ℓ_B and winning 11. Naturally, for the remaining 98 rounds reasonable player would decide to keep playing (exploit) with the "better" lever ℓ_B . This strategy might have two possible outcomes

- positive outcome: the lever ℓ_B is the machine ℓ_1 with $r_1 \sim \mathcal{N}(10, 3)$ and thus the mean reward would be $r = \frac{9+11}{2} + 10 \cdot 98 = 990\$$,
- negative outcome: the lever ℓ_B is the less beneficial lever ℓ_2 with $r_2 \sim \mathcal{N}(8, 4)$ giving mean reward $r = \frac{9+11}{2} + 8 \cdot 98 = 794\$$.

Since the player does not know whether the decision of choosing ℓ_B was correct or not, expected regret of this strategy would be $\mathbb{E}[\rho] = 1\,000 - [\frac{9+11}{2} + 0.5 \cdot (98 \cdot 10 + 98 \cdot 8)] = 1\,000 - 892 = 108$.

Explore-only strategy gave the player expected regret of 100, which is lower than exploit-only strategy with expected regret 108. However, we can see how sensitive the problem is. What if the player got lucky and in exploitation part chose the more winning machine? This dilemma is called *Exploration-Exploitation dilemma* and even though it is not the subject of this project, it is closely related to multi-armed bandit problem. It will come as no surprise that finding the balance between explore-exploit approaches is the key part of the most popular strategies for solving multi-armed bandit problem.

Strategies

In the example above were mentioned two basic strategies - **Exploration** and **Exploitation strategy**. In this section, some of the other strategies or solutions to multi-armed bandit problem are presented.

ε -first strategy

ε -first strategy extends the exploring part in exploitation strategy, which in general increases chances for choosing the more beneficial lever. Let $\varepsilon \in (0, 1)$.

The idea of ε -first strategy is to

1. explore (randomly choose levers) over the first εH rounds,
2. observe mean reward for each lever after first εH rounds played,
3. choose lever with higher mean reward,
4. exploit the chosen lever for the remaining $(1 - \varepsilon)H$ rounds.

The advantage of this method is without a doubt extended exploration part, which provides more clues about the hidden distributions F_i of each lever. On the other hand, since we know so little about the distributions of each lever, we can not with certainty determine the optimal range of exploration part. Setting the value of ε so we do not perform unnecessary exploration rounds on the expenses of exploitation rounds and the other way around is a difficult task to do and strongly depends on the distributions.

ε -greedy strategy

Another way to upgrade exploration/exploitation strategy is to use $\varepsilon \in (0, 1)$ as a threshold for re-decision. In ε -greedy strategy we continuously monitor our decision during the whole game. Firstly, lever is randomly chosen to be played for the first εH rounds. In $(\varepsilon H + 1)$ round the other lever is played. At this point, re-decision is made. If the reward obtained from other lever is lower than mean reward obtained during εH rounds, player continues to play with the same lever as in the εH rounds. In the opposite case, player randomly choose which lever is to be played in the following εH rounds.

In ε -greedy strategy is

1. randomly selected one of the levers,
2. exploited for εH rounds and calculated mean reward \tilde{r} after εH rounds,
3. in $(\varepsilon H + 1)$ round explored the second lever, obtaining levers reward r ,
 - if $\tilde{r} \geq r \rightarrow$ keep same lever as in exploitation $\rightarrow 2.$,
 - if $\tilde{r} < r \rightarrow$ randomly select new lever $\rightarrow 2.$ with new chosen lever,
4. the game ends after H rounds are played.

The contribution of this method to the improvement of previous strategies lies in periodic update of our decision after εH rounds. But again, in disadvantage, the series of playing the more beneficial lever could be accidentally interrupted by sporadic good win of the other lever in exploring round.

ε -decay strategy*****

ε -decay strategy is a special case of so called **ε -decreasing strategy**, where the $\varepsilon \in (0, 1)$ is not fixed, but it is a decreasing series of $(\varepsilon_i)_{i=1}^m$, where $\varepsilon_1 > \varepsilon_2 > \dots > \varepsilon_m$. In the ε -decay strategy, the lever of highest estimated mean is always pulled except when the least-taken lever is pulled with a probability of $\frac{4}{4+\beta t}$, where t is the number of times the least-taken lever has already been pulled.

LeastTaken strategy

1. randomly selected one of the levers,
2. exploited for $\varepsilon_i H$ rounds and calculated mean reward \tilde{r} after $\varepsilon_i H$ rounds,
3. in $(\varepsilon_i H + 1)$ the same lever is played with probability of $\frac{4}{4+m^2}$,
4. set ε_i as new ε_{i+1} ($\varepsilon_i > \varepsilon_{i+1}$) and return to point 2.,
5. the game ends after H rounds are played.

This method could be again upgraded by using more specific probability distribution for pulling the lever. Exaple of such method is SoftMax strategy:

SoftMax strategy

The SoftMax strategy consists of a random choice according to a Gibbs distribution. The lever k is chosen with probability

$$p_k = \frac{e^{\frac{\hat{\mu}_k}{\tau}}}{\sum_{i=1}^n e^{\frac{\hat{\mu}_i}{\tau}}},$$

where $\hat{\mu}_i$ is the estimated mean of the rewards brought by the lever i and $\tau \in \mathbb{R}^+$ is a parameter called the temperature. The choice of τ 's value is left to the user. More generally, all methods that choose levers according to a probability distribution reflecting how likely the levers are to be optimal, are called probability matching methods.

Interval Estimation strategy

A totally different approach to the exploration problem is to attribute to each lever an “optimistic reward estimate” within a certain confidence interval and to greedily choose the lever with the highest optimistic mean. Unobserved or infrequently observed levers will have an over-valued reward mean that will lead to further exploration of those levers. The more a lever is pulled and the closer its optimistic reward estimate will be to the true reward mean. This approach called Interval Estimation. To each lever is associated the $100 \cdot (1 - \alpha)\%$ reward mean upper bound where $\alpha \in (0, 1)$ is a parameter whose exact value is left to the user. At each round,

the lever of highest reward mean upper bound is chosen. Note that smaller α values lead to more exploration.

References

1. **Ritvik Kharkar**, A.A. [ritvikmath], Multi-Armed Bandit: Data Science Concepts (2020/09/23), YouTube, https://www.youtube.com/watch?v=e3L4VocZnnQ&ab_channel=ritvikmath
2. **Robert C. Gray**, A.A. [Academic Gamer], Multi-Armed Bandits: A Cartoon Introduction - DCBA #1 (2020/08/08), YouTube, https://www.youtube.com/watch?v=bkw6hWvh_3k&ab_channel=AcademicGamer
3. **Aleksandrs Slivkins**, Introduction to Multi-Armed Bandits, Microsoft Research NYC, Version September 2019, 2017-2019
4. **João Gama, Rui Camacho, Pavel B. Brazdil, Alípio Mário Jorge, Luís Torgo**, Machine Learning: ECML 2005, 16th European Conference on Machine Learning, Proceedings, Porto, Portugal, October 3-7, 2005, 0302-9743