

# Corporate Social Responsibility via Multi-Armed Bandits

Nwemadji Tiako Arsene Gibbs

Reinforcement Learning (RL) 2021-2022, QLS Diploma, ICTP.

Prof. Dr. Celani and Dr. Panizon

August 1, 2022

# Introduction

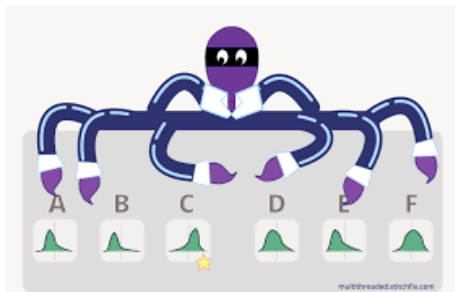
Corporate Social Responsibility (CSR), is a form of self-regulation that reflects a business's accountability and commitment to contributing to environmental and social measures.



Figure: Corporate Social Responsibility.

# Multi-Armed Bandit

Multi-armed Bandits.



# Setup and Notation

- ①  $K$ - arms
- ②  $T$ - number of time steps
- ③  $r_t$  reward obtained in round  $t$
- ④  $N_{i,t}$  number of times arm  $i$  is pull after  $t$  round
- ⑤  $\mu$  reward vector,  $\mu^*$  the expected reward of the optimal arm.
- ⑥  $\Delta_i = \mu^* - \mu_i$ , reward gap between an arm  $i$  and the optimal arm.
- ⑦  $f : [0, 1]^K \rightarrow [0, 1]^K$ , fairness function.
- ⑧  $\lambda$  transfer cost.

# Requirements and Examples of the Fairness Function

Requirements:

- ① Bounded  $\sum_{i=1}^K f(\mu)_i \leq 1$ .
- ② Lipschitz -  $\exists L > 0$  s.t.  $\forall \mu', \mu, \|f(\mu) - f(\mu')\|_1 \leq L \|\mu' - \mu\|_1$

Examples:

- Uniform  $f^{\text{uni}}(\mu)_i := \frac{1}{K}$
- Linear  $f^{\text{lin}}(\mu)_i := \frac{\mu_i}{K}$
- Step Uniform  $f^{\text{stp}}(\mu, d)_i := \frac{\theta(\mu_i - d)}{K}$
- Step Linear  $f^{\text{stp}}(\mu, d)_i := \mu_i \frac{\theta(\mu_i - d)}{K}$
- Softmax  $f^{\text{sft}}(\mu, c)_i := \frac{e^{c\mu_i}}{\sum_{j=1}^K e^{c\mu_j}}$ .

# Utility Function

With the help of the previous setting and given an algorithm  $ALG$ , the utility of the decision-maker is formally written as follow:

$$\mathcal{U}_{\lambda,f}(ALG; T) := \sum_{t=1}^T r_{i_t} - \lambda \sum_{i=1}^k \max\{Tf(\mu)_i - N_{i,T}, 0\}. \quad (1)$$

The performance of the algorithm is given by:

$$\mathcal{R}_{\lambda,f}(ALG; T) := \mathbb{E}(\mathcal{U}_{\lambda,f}(OPT; T)) - \mathbb{E}(\mathcal{U}_{\lambda,f}(ALG; T)), \quad (2)$$

where  $OPT$  is an algorithm maximizing the utility  $\mathcal{U}_{\lambda,f}(OPT; T)$ .

**Lemma 1.** Fix an arbitrary instance  $\langle K, T, \mu, f, \lambda \rangle$  and let  $OPT$  be an optimal algorithm for that instance. For every sub-optimal arm  $i$ , if  $\Delta_i < \lambda$  then  $OPT$  pulls  $i$  exactly  $Tf(\mu)_i$  times; if  $\Delta_i > \lambda$   $OPT$  does not pull  $i$  at all. If  $\Delta_i = \lambda$ ,  $OPT$  pulls arm  $i$  between zero and  $Tf(\mu)_i$  times.

**Theorem 1.** Fix any arbitrary instance  $\langle K, T, \mu, f, \lambda \rangle$  of R-O MAB, and let  $N = 8L^{2/3}T^{2/3}\log^{1/3}T$ . Algorithm 1 has a regret of  $O(KL^{2/3}T^{2/3}\log^{1/3}T)$ .

## Algorithm1: Fairness-Aware-ETC

```

0- Input:  $N$  - number of exploration rounds
1- for  $i = 1, \dots, K$  do
2-   pull arm  $i$  for  $N$  rounds
3- for  $i = 1, \dots, K$  do
4-   if  $\hat{\Delta}_i < \lambda$  do
5-     pull arm  $i$  for  $\max\{Tf(\hat{\mu}')_i - N, 0\}$  rounds
6- pull an arbitrary arm from  $\arg \max_{i \in [K]} \hat{\mu}$  until the execution ends.
  
```



## Example

Example: Consider a multi-armed bandit for which  $\mu = (1, \frac{1}{2}, \frac{1}{3})$ ,  $\lambda = 0.7$  and  $f(\mu)_i = \frac{1}{3}$ .

The classical ETC algorithm after discover that arm 1 is better will give:

$$T - 0.7T\left(\frac{1}{3} + \frac{1}{3}\right) = 0.533T.$$

While the optimal algorithm 1 after estimating the expected reward will gives:

$$\frac{T}{3} \cdot 1 + \frac{T}{3} \cdot \frac{1}{2} + \frac{T}{3} \cdot \frac{1}{3} = \frac{11}{18}T = 0.611T.$$

In this case the regret is:

$$R_{0.7} = 0.077T.$$

# Algorithm 2

**Theorem 2.** Fix any arbitrary instance  $\langle K, T, \mu, f, \lambda \rangle$  of R-O MAB, and let  $\alpha = K^{2/3} L^{2/3} T^{-1/3} \log^{1/3} T$ ,  $\beta = T^{-1/3} \log^{1/3} T$ . Then Algorithm 2 has a regret of  $O(K^{5/3} L^{2/3} T^{2/3} \log^{1/3} T)$ .

## Algorithm2: Self-regulated Utility Maximization

0- Input: Black-box bandit algorithm ALG, allowed approximation error parameter  $\alpha$  and  $\beta$ .

1-  $t = 1$

2- Initialize arms' data  $-N_i = 0, LCB(\Delta_i) = 0, UCB(\Delta_i) = 1$  for all  $i \in [K]$

3-  $C_1 = [0, 1]^K$  // Hyper-cube of  $\mu$  values

4- **While**  $\exists i \in [K]$  s.t.  $UCB(\Delta_i) > \lambda + \beta$  and  $LCB(\Delta_i) < \lambda - \beta$  **do** // Phase 1

5- Pull all arms once, update  $t$ , counters, confidence bounds and  $C_t$

6- **While**  $\exists i \in [K]$  s.t.  $\max_{\mu' \in C_t} f(\mu')_i - \min_{\mu' \in C_t} f(\mu')_i > \alpha$  and  $LCB(\Delta_i) < \lambda$  **do** // Phase 2:

7- Pull all arms once, update  $t$ , counters, confidence bounds and  $C_t$

8- **While**  $\exists i \in [K]$  s.t.  $LCB(\Delta_i) < \lambda$  and  $t < T$  **do** // Phase 3

9- Pull all arm  $i$  the minimal number of times so  $N_i \geq Tf(\hat{\mu}')_i$ , update  $t$  and counters.

10- Invoke ALG for the remaining rounds // Phase 4

# Setup of the Simulation

- 6 arms,  $\mu_i = 0.2 + (i - 1) \times 0.15$
- Transfer costs: 0, 0.4, 0.8
- Fairness function  $f^{sft}$
- Time Horizon: 10K, 50K, 100K and 200K
- Each experiment was repeated 50 times
- In case phase 1 or phase 2 had been made the black box algorithm is an algorithm which pull the arm with the highest expected reward
- If none exploration had been done when phase 4 start, Epsilon greedy algorithm was used as black-box algorithm.

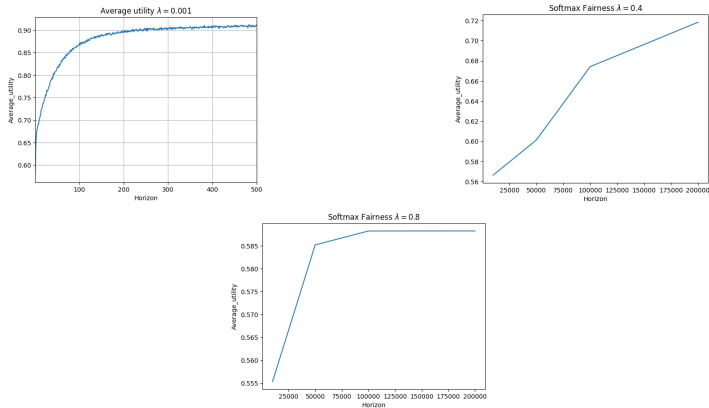


Figure: Average Utility for different transfer costs respect to the time horizon.

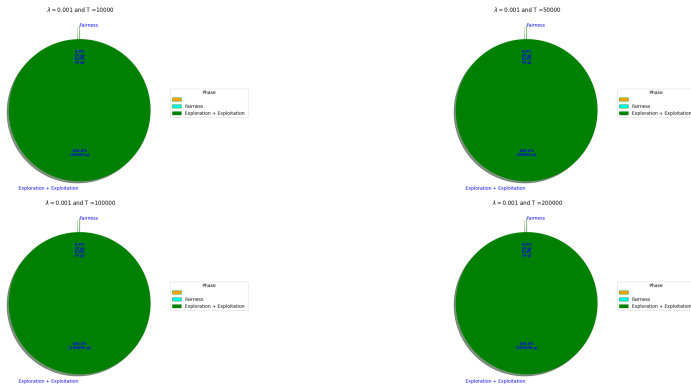


Figure: Round distribution per phase with  $\lambda = 0.001$  for softmax.

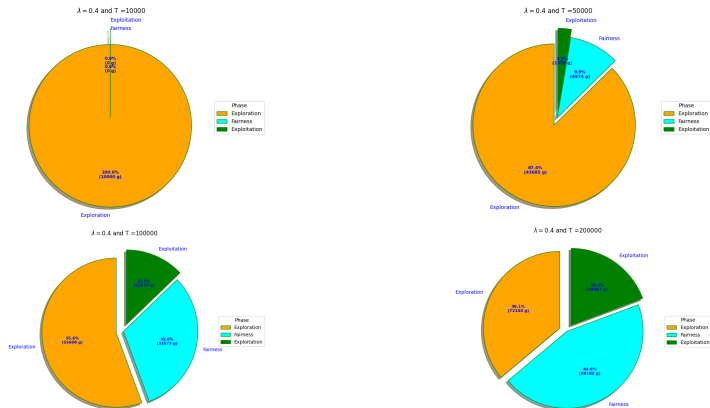


Figure: Round distribution per phase with  $\lambda = 0.4$  for softmax.

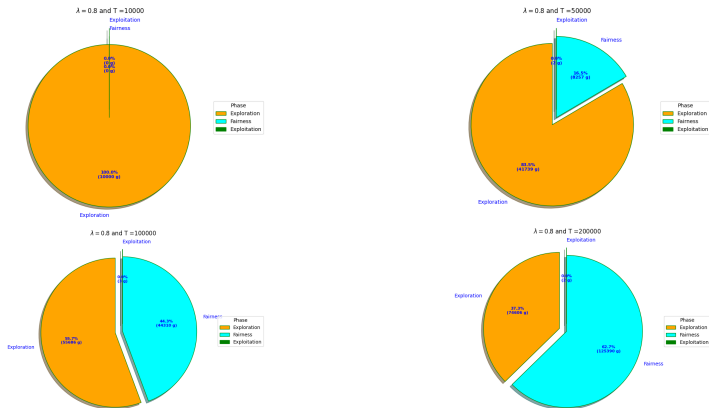


Figure: Round distribution per phase with  $\lambda = 0.8$  for softmax.



## Conclusion:

My goal was to recover the numerical result obtained on the paper for corporate social responsibility via multi-armed bandit by Rom, Ben-Porat et Sharit 2021. Despite the fact that I did not use the exact parameter they used on the paper, I succeeded to obtain most of their results.

# Question and Discussion

QUESTION