



POLITECNICO
MILANO 1863

FACTORED-REWARD BANDITS
WITH INTERMEDIATE OBSERVATIONS

Marco Mussi* Simone Drago* Marcello Restelli Alberto Maria Metelli
Politecnico di Milano

International Conference on Machine Learning @ Vienna, Austria

In the customary **Multi-Armed Bandit** framework, we consider a problem where:

- We have K **arms**, each representing an action
- The actions are **independent**
- There is **no structure** in the reward

However, in several cases, we may have:

- A structure in the **actions** and/or in the **reward** model
- Access to **intermediate effects** which may help the **learning process**

We consider the scenario in which we want to **sell a product online**:

- We have to choose a **price-budget pair**:
 - the **price** we set determines the users' propensity to buy (the so-called **conversion rate**)
 - the **advertising budget** we invest influences the number of potential customers that will be exposed (i.e., the number of **impressions**)
- We have access to **intermediate observations**:
 - the conversion rate, which depends on the price
 - the expected number of impressions, which depends on the budget
- Our objective is to maximize the revenue (i.e., **reward**) that is a function of the **product** between **intermediate observations**

- We can solve this problem using **standard Multi-Armed Bandit** techniques considering **price-budget couples** as actions
- However, if we look just at the reward and disregard this **factored** structure, the learning problem will:
 - present an **unnecessarily large action space**, including all the possible combinations of action components
 - suffer a **possibly amplified effect of the noise** in the reward due to the product of the noisy intermediate observations

- At every round $t \in \llbracket T \rrbracket$, we choose an **action vector**:

$$\mathbf{a}(t) = (a_1(t), \dots, a_d(t)) \in \mathcal{A} := \llbracket k_1 \rrbracket \times \dots \times \llbracket k_d \rrbracket$$

- $\forall i \in \llbracket d \rrbracket$ we have k_i options
 - d is the action vector dimension
- We observe a vector of d **intermediate observations** $\mathbf{x}(t) = (x_1(t), \dots, x_d(t))$ and receive as reward the **product of the observations** $r(t) = \prod_{i \in \llbracket d \rrbracket} x_i(t)$
 - The i^{th} component $x_i(t)$ of the intermediate observation vector $\mathbf{x}(t)$ is the effect of the i^{th} action component $a_i(t)$ in the action vector: $x_i(t) = \mu_{i, a_i(t)} + \epsilon_i(t)$
 - $\mu_{i, a_i(t)} \in [0, 1]$ is the **expected observation** of the i^{th} component $a_i(t)$
 - $\epsilon_i(t)$ is σ^2 -subgaussian noise

- An **optimal action vector** is:

$$\mathbf{a}^* = (a_1^*, \dots, a_d^*) \in \arg \max_{\mathbf{a}=(a_1, \dots, a_d) \in \mathcal{A}} \prod_{i \in [d]} \mu_{i, a_i}$$

and we abbreviate $\mu_i^* = \mu_{i, a_i^*}, \forall i \in [d]$

- We define the **suboptimality gaps** related to:

- the i^{th} **action component** $\Delta_{i, a_i} := \mu_i^* - \mu_{i, a_i}$ for $a_i \in [k_i]$
- the **action vector** $\mathbf{a} = (a_1, \dots, a_d) \in \mathcal{A}$ as $\Delta_{\mathbf{a}} := \prod_{i \in [d]} \mu_i^* - \prod_{i \in [d]} \mu_{i, a_i}$

- The **goal** of an algorithm \mathcal{A} is to minimize the **expected cumulative regret**:

$$\mathbb{E}[R_T(\mathcal{A}, \underline{\nu})] := \mathbb{E} \left[T \prod_{i \in [d]} \mu_i^* - \sum_{t \in [T]} \prod_{i \in [d]} \mu_{i, a_i(t)} \right] = \mathbb{E} \left[\sum_{t \in [T]} \Delta_{\mathbf{a}(t)} \right]$$

Theorem (Worst-Case Lower Bound)

For every algorithm \mathfrak{A} , there exists an FRB $\underline{\nu}$ such that for $T \geq \mathcal{O}(d^2)$, \mathfrak{A} suffers an expected cumulative regret of at least:

$$\mathbb{E} [R_T(\mathfrak{A}, \underline{\nu})] \geq \frac{\sigma}{4\sqrt{2}} \sum_{i \in \llbracket d \rrbracket} \sqrt{k_i T}.$$

In particular, if $k_i =: k$ for every $i \in \llbracket d \rrbracket$, we have:

$$\mathbb{E} [R_T(\mathfrak{A}, \underline{\nu})] \geq \Omega(\sigma d \sqrt{kT}).$$

Theorem (Instance-Dependent Lower Bound)

For every consistent algorithm \mathfrak{A} and FRB $\underline{\nu}$ with unique optimal arm $\mathbf{a}^* \in \mathcal{A}$ it holds:

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}[R_T(\mathfrak{A}, \underline{\nu})]}{\log T} \geq \underline{C}(\underline{\nu}) = \min_{(L_{\mathbf{a}})_{\mathbf{a} \in \mathcal{A} \setminus \{\mathbf{a}^*\}}} \sum_{\mathbf{a} \in \mathcal{A} \setminus \{\mathbf{a}^*\}} L_{\mathbf{a}} \Delta_{\mathbf{a}}$$

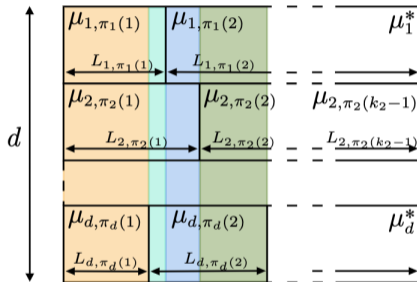
s.t. $L_{i,j} = \sum_{\mathbf{a} \in \mathcal{A} \setminus \{\mathbf{a}^*\}, a_i=j} L_{\mathbf{a}}, \quad \forall i \in \llbracket d \rrbracket, j \in \llbracket k_i \rrbracket \setminus \{a_i^*\}$

$$L_{i,j} \geq \frac{2\sigma^2}{\Delta_{i,j}^2}, \quad \forall i \in \llbracket d \rrbracket, j \in \llbracket k_i \rrbracket \setminus \{a_i^*\}$$

$$L_{\mathbf{a}} \geq 0, \quad \forall \mathbf{a} \in \mathcal{A} \setminus \{\mathbf{a}^*\}.$$

- We consider $L_{i,j} = \mathbb{E}[N_{i,j}]/\log T$ to handle the asymptotic nature of the bound

- To solve the optimization problem, we have to search for the **best way** to arrange the pulls
- We can make use of **rearrangement inequality** for integrals to find the best solution (Luttinger and Friedberg, 1976)



- We present **Factored Upper Confidence Bound (F-UCB)**
- F-UCB performs a **UCB-like** exploration (Auer et al., 2002) **independently** for every dimension $i \in \llbracket d \rrbracket$
- Then, we study its **theoretical guarantees**

Algorithm: F-UCB.

Input : Exploration Parameter α , Subgaussian proxy σ , Action component size k_i , $\forall i \in \llbracket d \rrbracket$

1 Initialize $N_{i,a_i}(0) \leftarrow 0$, $\hat{\mu}_{i,a_i}(0) \leftarrow 0 \quad \forall a_i \in \llbracket k_i \rrbracket$, $i \in \llbracket d \rrbracket$

2 **for** $t \in \llbracket T \rrbracket$ **do**

3 Select $\mathbf{a}(t) \in \arg \max_{\mathbf{a}=(a_1, \dots, a_d)^T \in \mathcal{A}} \prod_{i \in \llbracket d \rrbracket} \text{UCB}_{i,a_i}(t)$ where $\text{UCB}_{i,a_i}(t) = \hat{\mu}_{i,a_i}(t-1) + \sigma \sqrt{\frac{\alpha \log t}{N_{i,a_i}(t-1)}}$

4 Play $\mathbf{a}(t)$ and observe $\mathbf{x}(t) = (x_1(t), \dots, x_d(t))^T$

5 Update $\hat{\mu}_{i,a_i}(t)$ and $N_{i,a_i}(t)$ for every $i \in \llbracket d \rrbracket$

6 **end**

Theorem (Worst-Case Upper Bound for F-UCB)

For any FRB $\underline{\nu}$, F-UCB with $\alpha > 2$ suffers an expected regret bounded as:

$$\mathbb{E} [R_T(\text{F-UCB}, \underline{\nu})] \leq 4\sigma \sum_{i \in [d]} \sqrt{\alpha k_i T \log T} + g(\alpha) \sum_{i \in [d]} k_i.$$

In particular, if $k_i =: k$, for every $i \in [d]$, we have:

$$\mathbb{E} [R_T(\text{F-UCB}, \underline{\nu})] \leq \tilde{O}(\sigma d \sqrt{kT}).$$

Theorem (Instance-Dependent Upper Bound for F-UCB)

For a given FRB $\underline{\nu}$, F-UCB with $\alpha > 2$ suffers an expected regret bounded as:

$$\mathbb{E}[R_T(\text{F-UCB}, \underline{\nu})] \leq \bar{C}(\text{F-UCB}, \underline{\nu}) = \max_{(N_{\mathbf{a}})_{\mathbf{a} \in \mathcal{A}}} \sum_{\mathbf{a} \in \mathcal{A} \setminus \{\mathbf{a}^*\}} N_{\mathbf{a}} \Delta_{\mathbf{a}}$$

$$\text{s.t. } N_{i,j} = \sum_{\mathbf{a} \in \mathcal{A} \setminus \{\mathbf{a}^*\}, a_i=j} N_{\mathbf{a}}, \quad \forall i \in \llbracket d \rrbracket, j \in \llbracket k_i \rrbracket \setminus \{a_i^*\}$$

$$N_{i,j} \leq \frac{4\alpha\sigma^2 \log T}{\Delta_{i,j}^2} + g(\alpha), \quad \forall i \in \llbracket d \rrbracket, j \in \llbracket k_i \rrbracket \setminus \{a_i^*\}$$

$$\sum_{\mathbf{a} \in \mathcal{A}} N_{\mathbf{a}} = T$$

$$N_{\mathbf{a}} \geq 0, \quad \forall \mathbf{a} \in \mathcal{A}$$

Corollary (Explicit Instance-Dependent Upper Bound for F-UCB)

For a given FRB $\underline{\nu}$, F-UCB with $\alpha > 2$ suffers an expected regret bounded by:

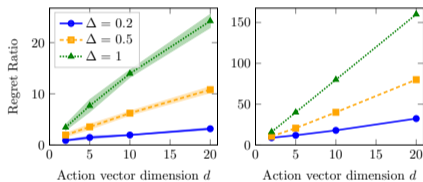
$$\begin{aligned} \mathbb{E} [R_T(F\text{-UCB}, \underline{\nu})] &\leq \overline{C}(F\text{-UCB}, \underline{\nu}) \\ &\leq 4\alpha\sigma^2 \log T \sum_{i \in [d]} \mu_{-i}^* \sum_{j \in [k_i] \setminus \{a_i^*\}} \Delta_{i,j}^{-1} + g(\alpha) \sum_{i \in [d]} k_i, \end{aligned}$$

where $\mu_{-i}^* = \prod_{l \in [d] \setminus \{i\}} \mu_l^* \leq 1$ for every $i \in [d]$.

- For $T \rightarrow +\infty$, we observe that:

$$\frac{\overline{C}(\text{F-UCB}, \underline{\nu})}{\underline{C}(\underline{\nu}) \log T} \leq \frac{2d\alpha\Delta}{1 - (1 - \Delta)^d} \stackrel{\Delta \rightarrow 1}{\approx} 2\alpha d$$

- F-UCB performs worse than the lower bound, with an **additional dependence on d**
- In the figure, we compare:
 - (left) the ratio between the regret obtained by running F-UCB and the instance-dependent lower bound
 - (right) the bound above



- **F-UCB** does not enjoy instance-dependent optimality due to the **lack of synchronization** over the components of the action vector
- To overcome this problem, we propose **F-Track**
- **F-Track** is an algorithm which computes and **tracks the lower bound** (Lattimore and Szepesvari, 2017)

Algorithm: F-Track.

Input : Warm-up sample size N_0 , Threshold ϵ_T , Action component size k_i , $\forall i \in \llbracket d \rrbracket$,

- 1 $t \leftarrow 1$
- 2 **while** $\min_{i \in \llbracket d \rrbracket} \min_{j \in \llbracket k_i \rrbracket} N_{i,j}(t) < N_0$ **do**
- 3 | Pull action vector $\mathbf{a}(t)$ with $a_i(t) = (t - 1) \bmod k_i + 1$ for all $i \in \llbracket d \rrbracket$, $t \leftarrow t + 1$
- 4 **end**
- 5 $T_{\text{warm-up}} \leftarrow t - 1$
- 6 Estimate the suboptimality gaps $\forall i \in \llbracket d \rrbracket, j \in \llbracket k_i \rrbracket$: $\hat{\Delta}_{i,j} := \max_{j' \in \llbracket k_i \rrbracket} \hat{\mu}_{i,j'}(T_{\text{warm-up}}) - \hat{\mu}_{i,j}(T_{\text{warm-up}})$
- 7 Compute the number of pulls $\hat{N}_{i,j} = 2\sigma^2 f_T(1/T) \hat{\Delta}_{i,j}^{-2}$ for every action component $i \in \llbracket d \rrbracket$ and $j \in \llbracket k_i \rrbracket$
- 8 Compute the number of pulls $\hat{N}_{\mathbf{a}}$ for every action vector $\mathbf{a} \in \mathcal{A}$ by solving the LP of the ID Lower Bound
- 9 **while** $t \leq T$ and $\max_{i \in \llbracket d \rrbracket, j \in \llbracket k_i \rrbracket} |\hat{\mu}_{i,j}(T_{\text{warm-up}}) - \hat{\mu}_{i,j}(t - 1)| \leq 2\epsilon_T$ **do**
- 10 | Pull action vector $\mathbf{a}(t) \in \arg \min \{N_{\mathbf{a}}(t) : \mathbf{a} \in \mathcal{A} \text{ and } N_{\mathbf{a}}(t) \leq \hat{N}_{\mathbf{a}}\}$, $t \leftarrow t + 1$
- 11 **end**
- 12 Discard all data and play F-UCB until $t = T$

Theorem (Instance-Dependent Upper Bound for F-Track)

For any FRB $\underline{\nu}$, F-Track run with:

$$f_T(\delta) := \left(1 + \frac{1}{\log T}\right) \left(c \log \log T + \log \left(\frac{1}{\delta}\right)\right),$$
$$N_0 = \lceil \sqrt{\log T} \rceil \quad \text{and} \quad \epsilon_T = \sqrt{\frac{2\sigma^2 f_T(1/\log T)}{N_0}},$$

suffers an expected regret of:

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}[R_T(\text{F-Track}, \underline{\nu})]}{\log T} = \underline{C}(\underline{\nu}).$$

- We presented the **Factored-Reward Bandits**, where we perform a **set of actions**, whose effects can be **observed**, and the reward is the **product** of those effects
- We characterized the **statistical complexity** of the setting from both the **worst-case** and **instance-dependent** perspectives
- We presented **F-UCB**, and we characterized its **instance-dependent** and **worst-case** guarantees and we discuss its **instance-dependent limitations**
- To overcome the F-UCB's limitations, we presented **F-Track**, which shows asymptotical **instance-dependent** optimality

- Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2):235–256.
- Lattimore, T. and Szepesvari, C. (2017). The end of optimism? an asymptotic analysis of finite-armed linear bandits. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 728–737. PMLR.
- Luttinger, J. and Friedberg, R. (1976). A new rearrangement inequality for multiple integrals. *Archive for Rational Mechanics and Analysis*, 61:45–64.