# FACTORED-REWARD BANDITS WITH INTERMEDIATE OBSERVATIONS

MARCO MUSSI*, SIMONE DRAGO*, MARCELLO RESTELLI AND ALBERTO MARIA METELLI

{marco.mussi, simone.drago, marcello.restelli, albertomaria.metelli}@polimi.it

POLITECNICO MILANO 1863 — RL³

## EXAMPLE: JOINT PRICING–ADVERTISING



Actions: Price, Budget → Intermediate Observations: Conversion Rate, Number of Impressions → × → Reward
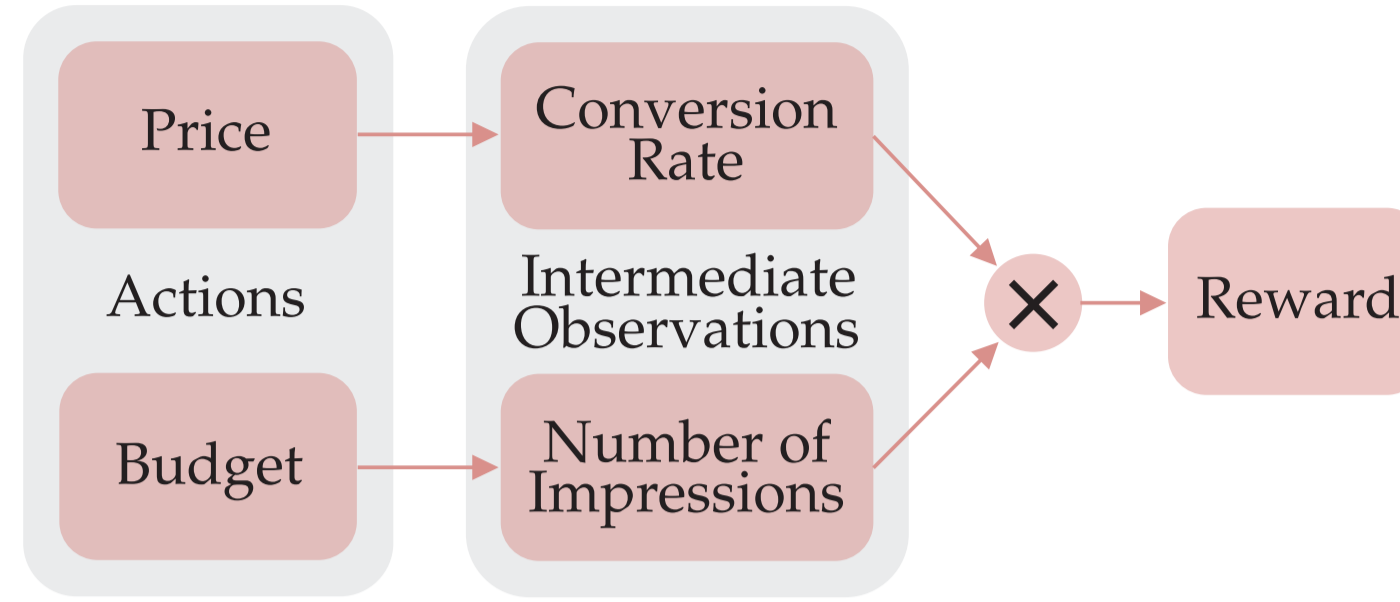
## WHY NOT STANDARD MAB?

We can solve this problem using **standard** Multi-Armed Bandit techniques considering the **price-budget couples** as actions, at the cost of an:

- **unnecessarily large action space** ($|\mathcal{A}| = \prod_{i \in \llbracket d \rrbracket} k_i$)
- **amplified heavy-tailed noise** effect

## FACTORED-REWARD BANDITS (FRB)

We choose an **action vector**:

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

We observe a vector of $d$ **intermediate observations**:

$$\mathbf{x}(t) = (x_1(t), \dots, x_d(t))$$

with:

$$x_i(t) = \underbrace{\mu_{i,a_i(t)}}_{\substack{\text{Expected intermediate} \\ \text{observation of } a_i(t) \\ (\text{with } \mu_{i,j} \in [0,1])}} + \underbrace{\epsilon_i(t)}_{\substack{\sigma^2\text{-subgaussian} \\ \text{noise}}}$$

We receive a **reward**: $r(t) = \prod_{i \in \llbracket d \rrbracket} x_i(t)$

We consider $k_i = k$, $\forall i \in \llbracket d \rrbracket$ for *simplicity*.

### LEARNING PROBLEM

Optimal **action vector**:

$\mathbf{a}^* = (a_1^*, \dots, a_d^*) \in \bigtimes_{i \in \llbracket d \rrbracket} \arg\max_{a_i \in \llbracket k_i \rrbracket} \mu_{i,a_i}$

Optimal **expected reward**:

$\prod_{i \in \llbracket d \rrbracket} \max_{a_i \in \llbracket k_i \rrbracket} \mu_{i,a_i} = \prod_{i \in \llbracket d \rrbracket} \mu_i^* = \mu^*$

Suboptimality gaps: $\Delta_{i,a_i} := \mu_i^* - \mu_{i,a_i}$

Goal is to minimize the **expected cumulative regret**:

$\mathbb{E}[R_T(\mathfrak{A}, \underline{\boldsymbol{\nu}})] = T\mu^* - \mathbb{E}\left[\sum_{t \in \llbracket T \rrbracket} \prod_{i \in \llbracket d \rrbracket} \mu_{i,a_i(t)}\right]$

## LOWER BOUNDS

### WORST-CASE LOWER BOUND

$$\mathbb{E}[R_T(\mathfrak{A}, \underline{\boldsymbol{\nu}})] \geqslant \Omega\left(\sigma d \sqrt{kT}\right)$$

### INSTANCE-DEPENDENT LOWER BOUND

$$\liminf_{T \to +\infty} \frac{\mathbb{E}[R_T(\mathfrak{A}, \underline{\boldsymbol{\nu}})]}{\log T} \geqslant \underline{C}(\underline{\boldsymbol{\nu}})$$

- Every algorithm $\mathfrak{A}$ has to pull at least:
  $\frac{\mathbb{E}[N_{i,j}]}{\log T} \geqslant \frac{2\sigma^2}{\Delta_{i,j}^2}$
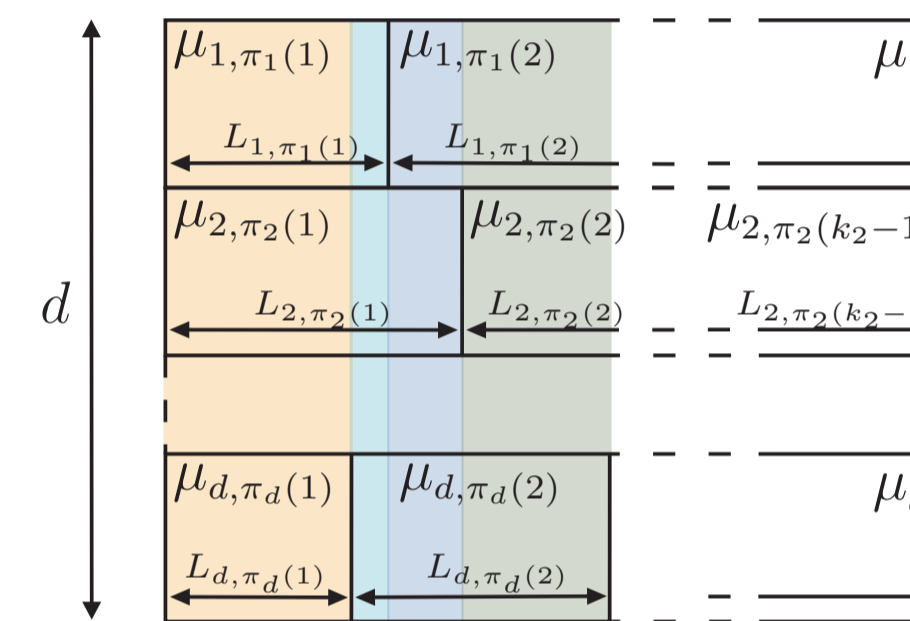  times every suboptimal action component
- The lower bound is obtained by finding the combination of pulls minimizing the regret
- The naïve approach is to solve a **Linear Programming optimization problem**

#### EFFICIENT SOLUTION TO THE LP

Using **Rearrangement Inequality**
$\mathcal{O}(dk \log(k))$ complexity



## SOLUTION 1: FACTORED UPPER CONFIDENCE BOUND (F-UCB)

F-UCB is an **anytime optimistic regret minimization** algorithm that plays over the $d$ different dimensions **independently**. In every dimension, the algorithm plays the action defined as:

$$\mathbf{a}(t) = \arg\max_{(a_1, \dots, a_d) \in \mathcal{A}} \prod_{i \in \llbracket d \rrbracket} \text{UCB}_{i,a_i}(t)$$

where the **optimistic index** is: $\text{UCB}_{i,a_i}(t) = \hat{\mu}_{i,a_i}(t-1) + \sigma \sqrt{\frac{\alpha \log t}{N_{i,a_i}(t-1)}}$

### WORST-CASE UPPER BOUND

$$\mathbb{E}[R_T(\text{F-UCB}, \underline{\boldsymbol{\nu}})] \leqslant \tilde{\mathcal{O}}\left(\sigma d \sqrt{kT}\right)$$

### INSTANCE-DEPENDENT UPPER BOUND

#### IMPLICIT UPPER BOUND

- F-UCB pulls at most:
  $\mathbb{E}[N_{i,j}] \leqslant \frac{4\alpha\sigma^2 \log T}{\Delta_{i,j}^2}$
  times every suboptimal action component
- We want to find the worst combination of pulls
- Again, the naïve approach is to solve a Linear Programming optimization problem
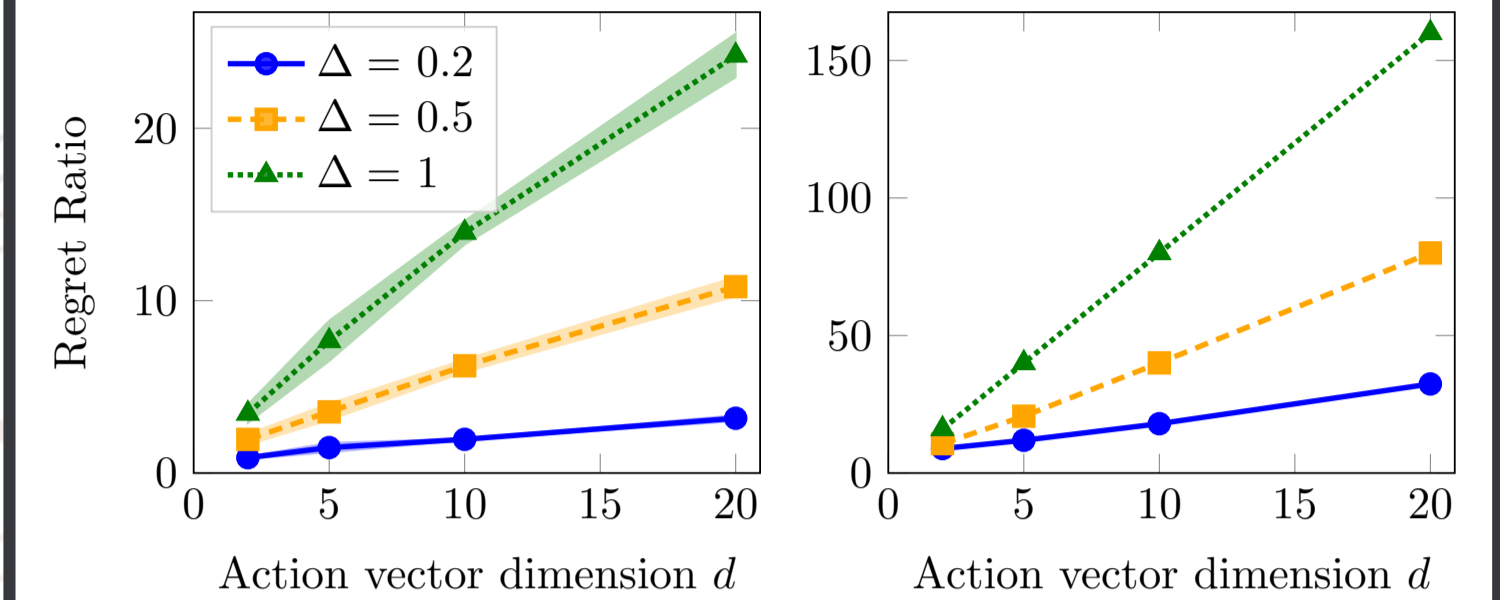
#### EXPLICIT UPPER BOUND
*(Rearrangement Inequality, opposite direction)*

$\mathbb{E}[R_T(\text{F-UCB}, \underline{\boldsymbol{\nu}})]$
$\leqslant 4\alpha\sigma^2 \log T \sum_{i \in \llbracket d \rrbracket} \mu_{-i}^* \sum_{j \in \llbracket k \rrbracket \setminus \{a_i^*\}} \Delta_{i,j}^{-1}$

where $\mu_{-i}^* = \prod_{l \in \llbracket d \rrbracket \setminus \{i\}} \mu_l^* \leqslant 1, \forall i \in \llbracket d \rrbracket$

## SOLUTION 2: F-TRACK

### F-UCB IS INSTANCE-DEPENDENT SUBOPTIMAL IN $d$



### SOLUTION: F-TRACK

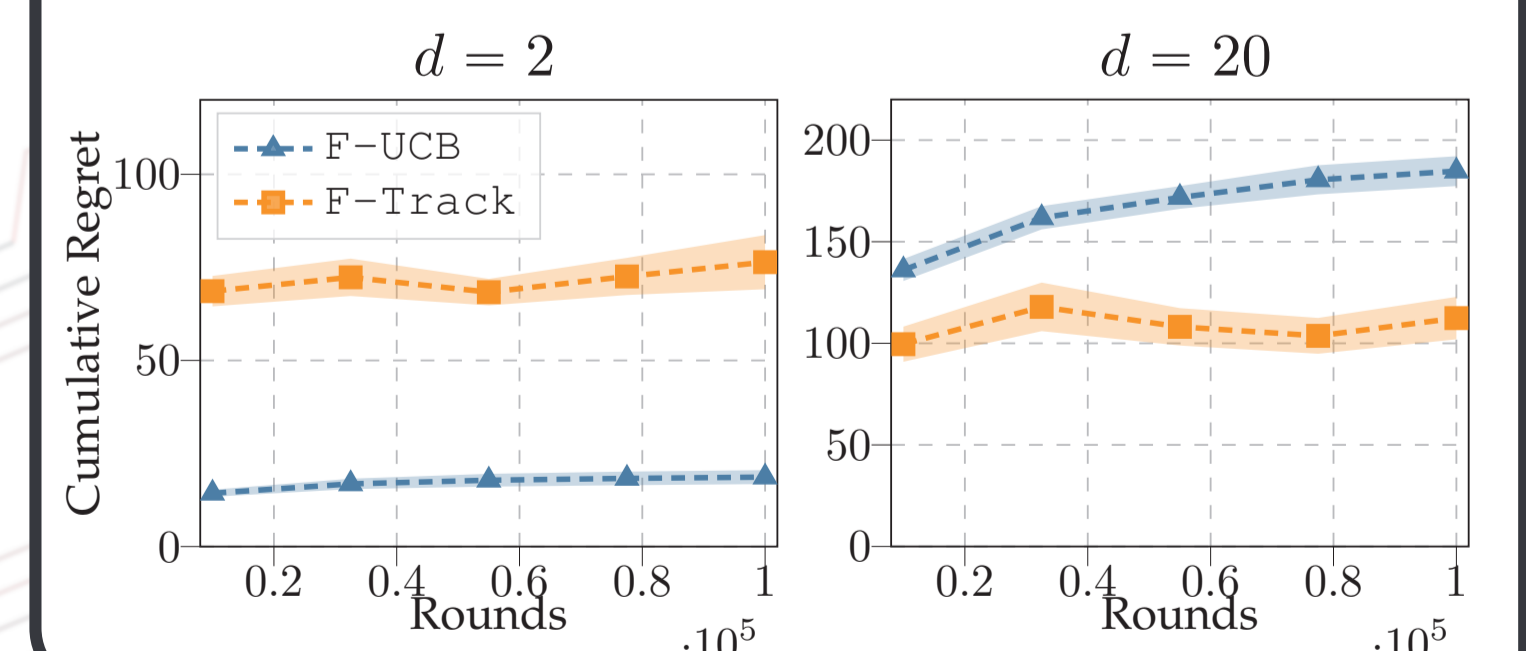F-Track **coordinates** among the $d$ dimensions in three phases:

1. **Warm-up**: Play action vectors in round robin until every action component has been pulled at least a minimum amount of times

2. **LB Matching**: Use warm-up data to compute estimates of $\hat{\mu}_{i,j}$ and $\hat{\Delta}_{i,j}$. Solve the lower bound LP to define a pull schedule

3. **Recovery**: If, during phase 2, the estimation error of any $\hat{\mu}_{i,j}$ is discovered to invalidate the scheduling, fall back to F-UCB until $t = T$

### INSTANCE-DEPENDENT UPPER BOUND

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

## EXPERIMENTAL RESULTS

Comparison between **F-UCB** and **F-Track** for different values of $d$. Setting: $k = 2$, $\mu^* = 1$, $\Delta = 0.7$.



## REFERENCES

S. Bubeck, N. Cesa-Bianchi, and G. Lugosi. Bandits with heavy tail. *IEEE Trans. Inf. Theory*, 2013.

T. Lattimore and C. Szepesvári. The end of optimism? an asymptotic analysis of finite-armed linear bandits. In *AISTATS*, 2017.

J. Zimmert and Y. Seldin. Factored bandits. In *NeurIPS*, 2018.