

ROTTING BANDITS ARE NOT HARDER THAN STOCHASTIC ONES

J. SEZNEC, A. LOCATELLI, A. CARPENTIER, A. LAZARIC, M. VALKO

julien.seznec@lelivrescolaire.fr

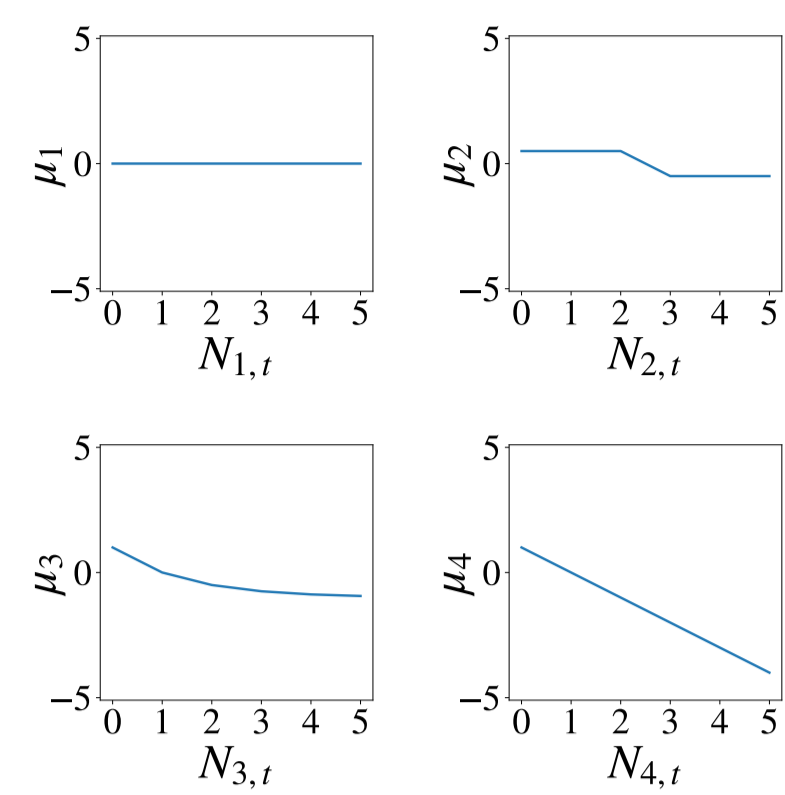
Setup



BANDITS:

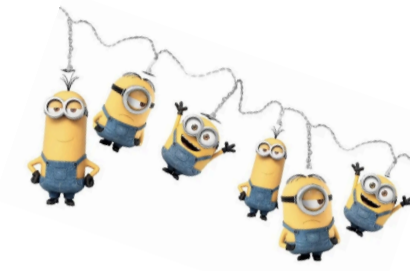
At each round t :

- SELECT an action i
- RECEIVE noisy reward



ROTTING HYPOTHESIS:

- Each time we pull an arm, reward decay
- Maximum decay between two pulls: L



GOAL: Maximize cumulative reward $\sum_{i \in \mathcal{K}} \sum_{n=0}^{N_{i,T}^* - 1} \mu_i(n)$

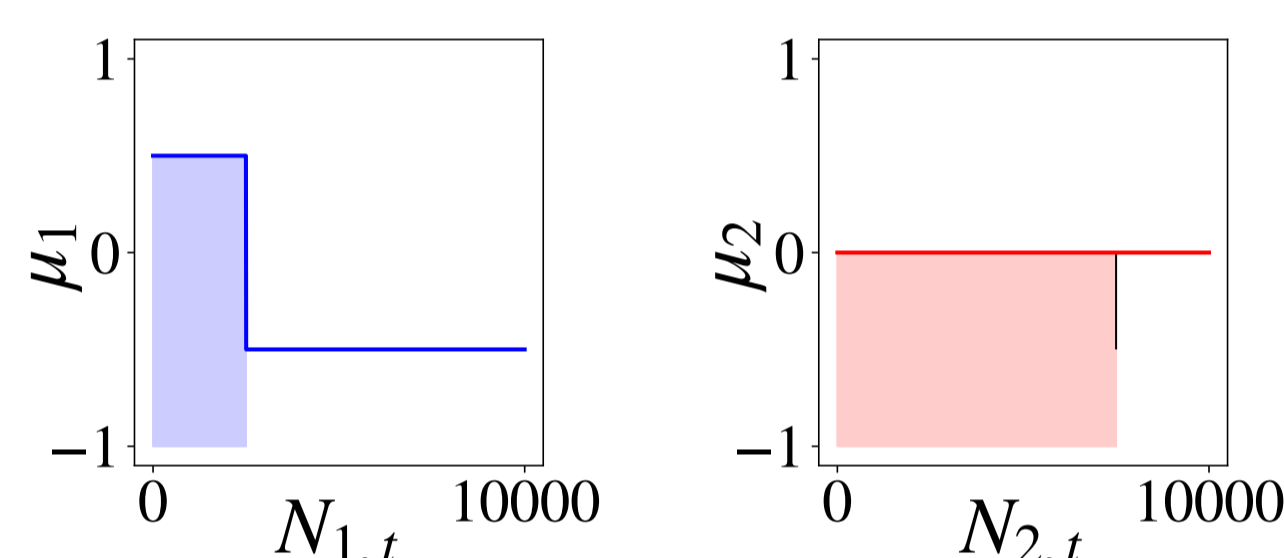
APPLICATIONS: education, economics ...

What's new?

FEWA algorithm

- ★ **Minimax optimal.** $T^{2/3} \rightarrow T^{1/2}$
- ★ **Solves an open problem**
- ★ **First problem dependent guarantee - recovers bandits**
- ★ **Adaptive to L - does not need it as an input**

Prior work



- Heidari et al. (2016)

- **Optimal oracle - knows μ_i**
→ Select the arm with largest next reward
Define regret against this optimal oracle

$$R_T(\pi) = \sum_{i \in \text{UP}} \sum_{s=N_{i,T}^*+1}^{N_{i,T}^*} \mu_i(s) - \sum_{i \in \text{OP}} \sum_{s=N_{i,T}^*+1}^{N_{i,T}^*} \mu_i(s)$$

- **No noise - but unknown μ_i**
→ Select the arm with largest last reward
→ Is minimax optimal.



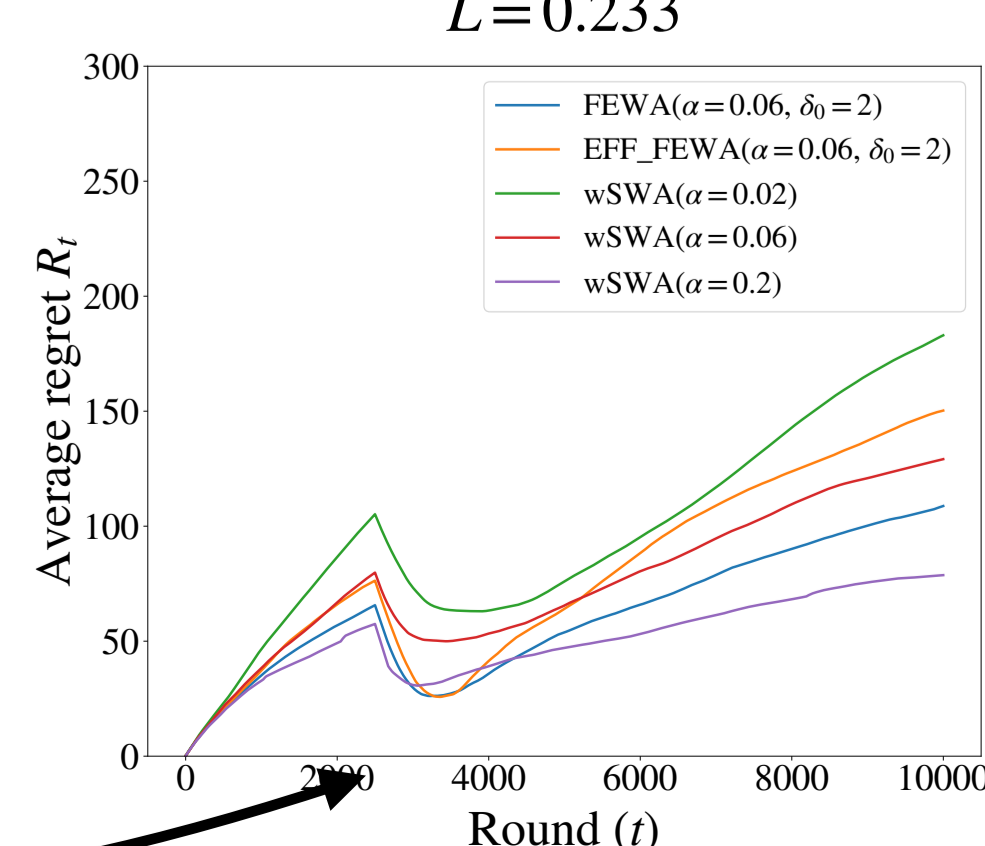
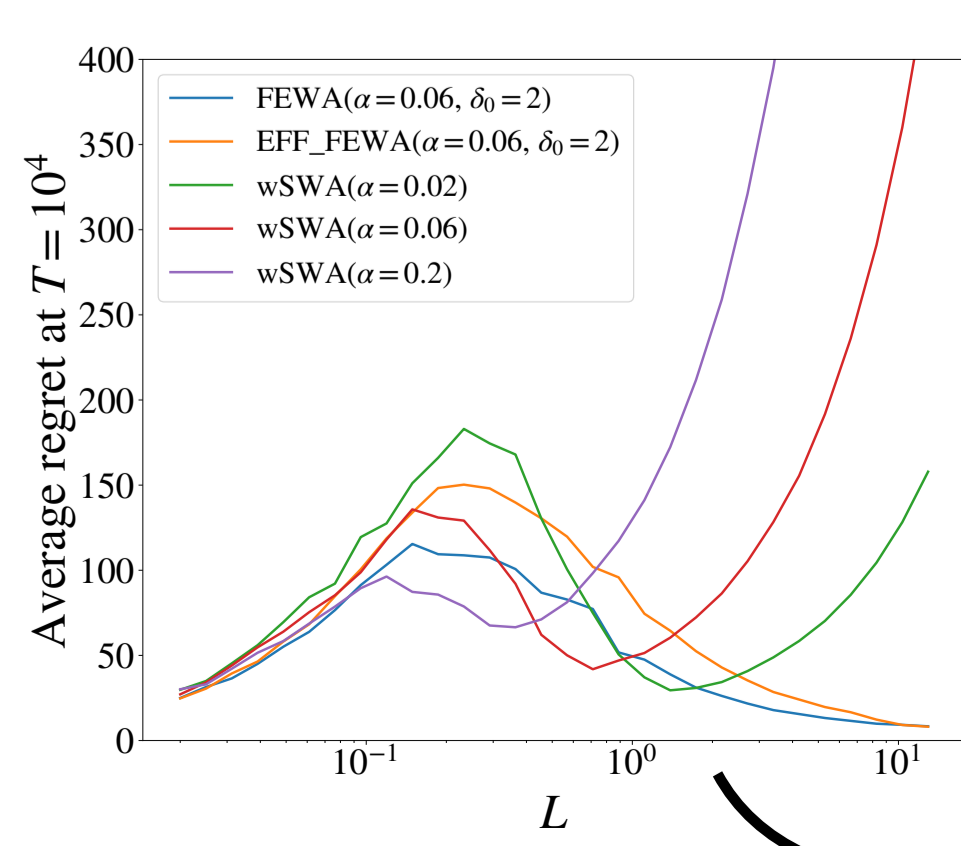
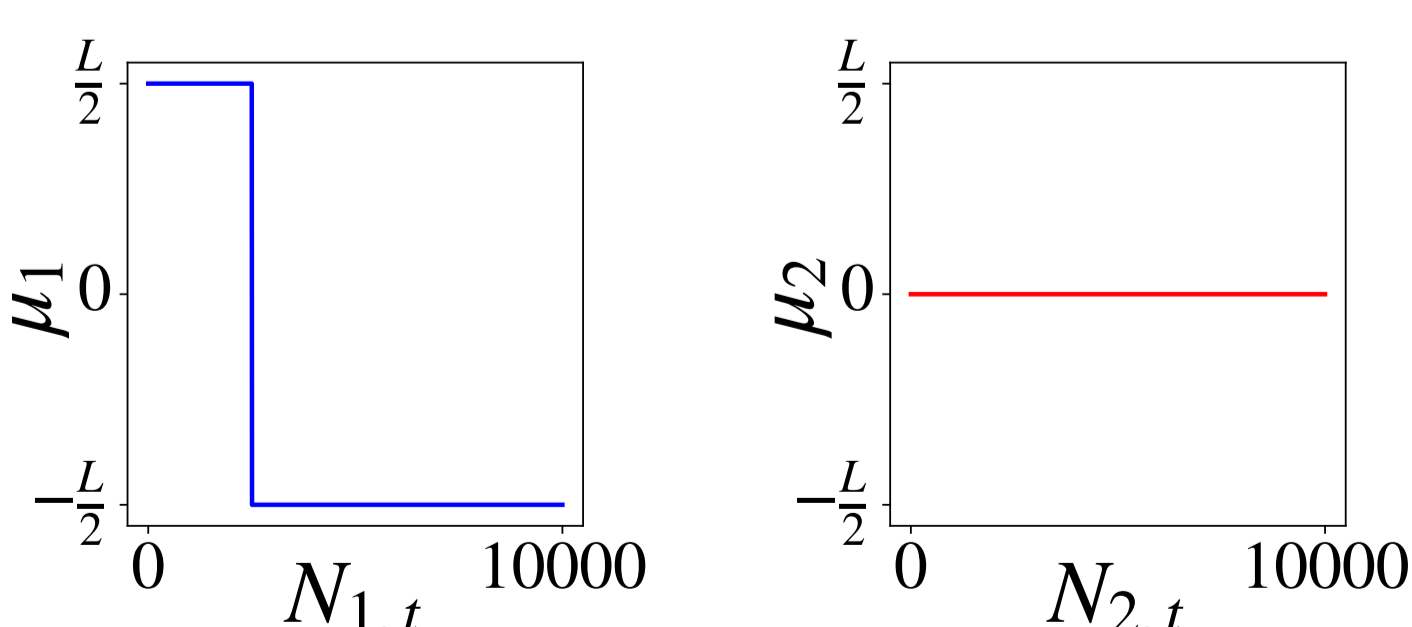
- Levine et al. (2017)

- Noisy reward *sliding window average*:
→ Select the arm with largest average of the h last reward sample $\mathbb{E}[R_T(G)] \leq KL$

- Optimize h for the bias-variance trade-off
 $\mathbb{E}[R_T(\pi_{\text{wSWA}})] = \tilde{O}(L^{1/3} \sigma^{2/3} K^{1/3} T^{2/3})$

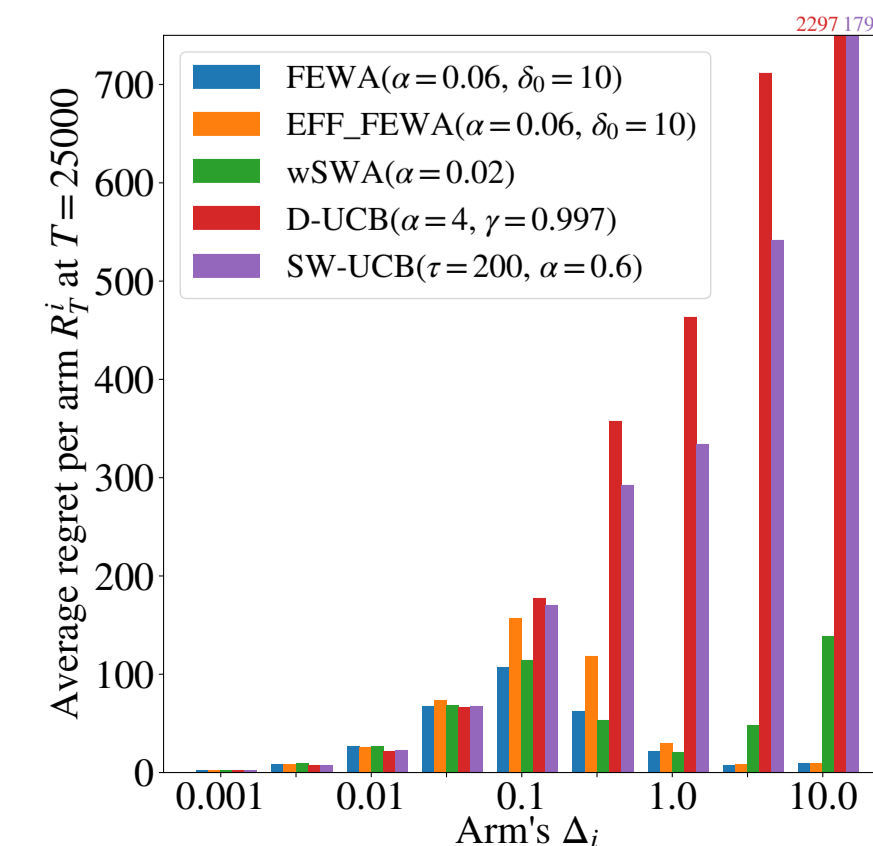
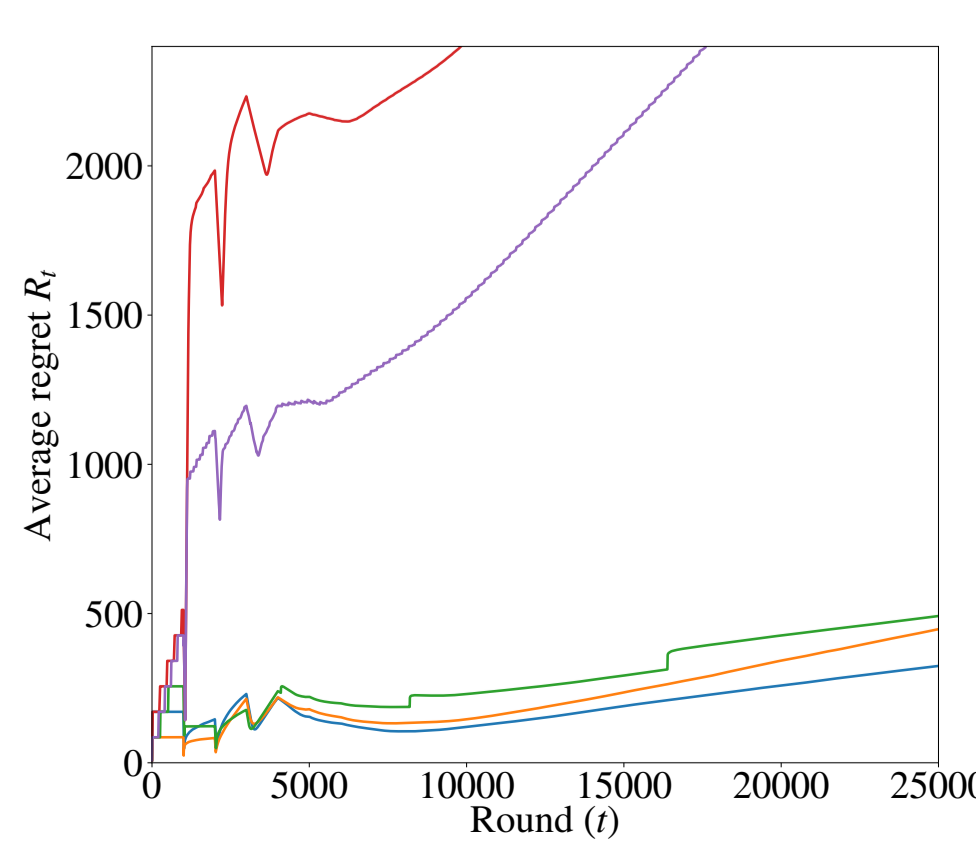
Experiments

2 arms - Minimal single drop experiment

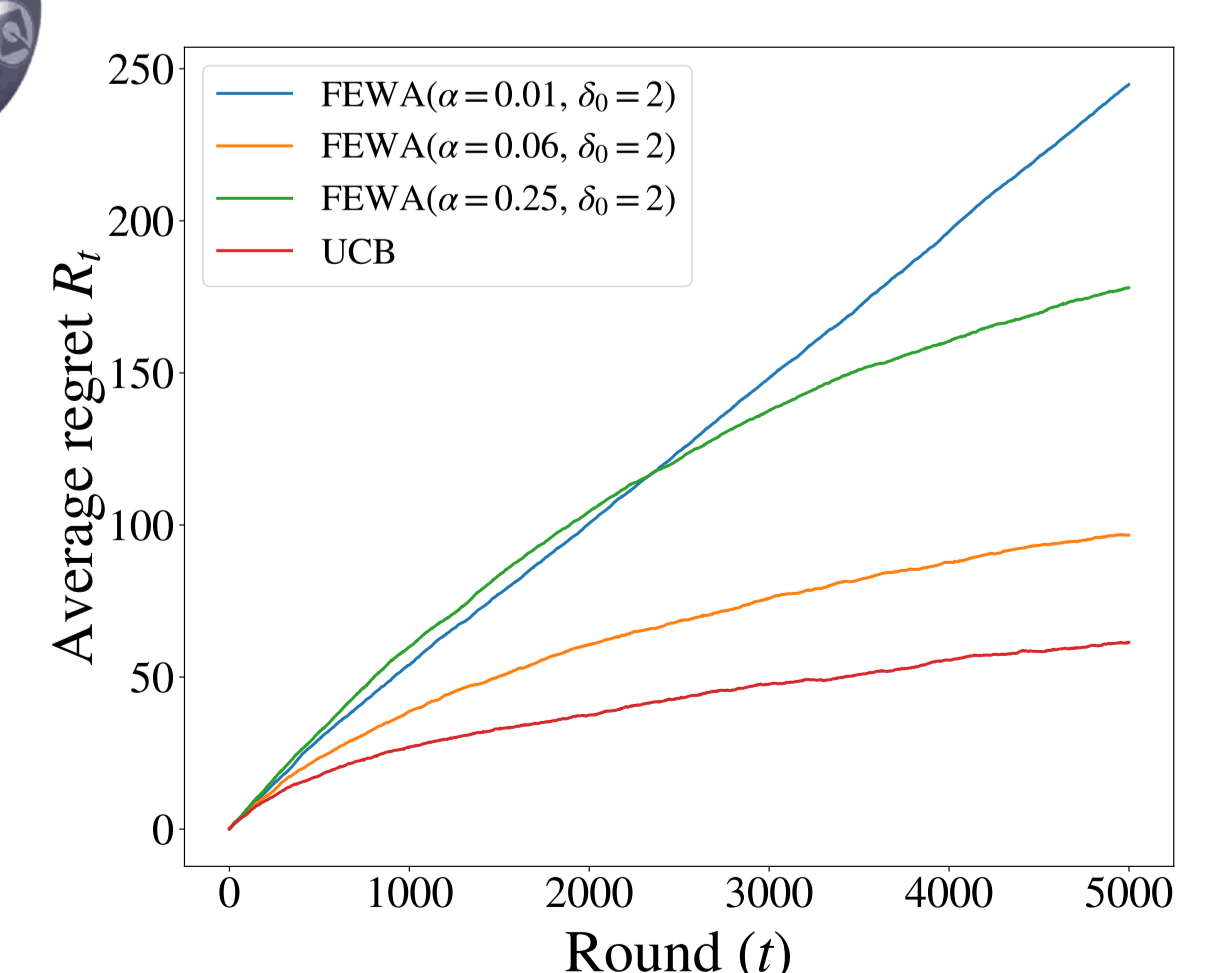


10 arms - Adapting to multiple decays

- 1 constant arm
- 9 arms with abrupt decay at 1000 pulls
- Geometric sequence of decays: 0.002 → 20



Competing against UCB1



- 1) Filtering policy < UCB index policy
- 2) More possible events → Looser CB

Algorithm

FEWA: Filtering on Expanding Window Average

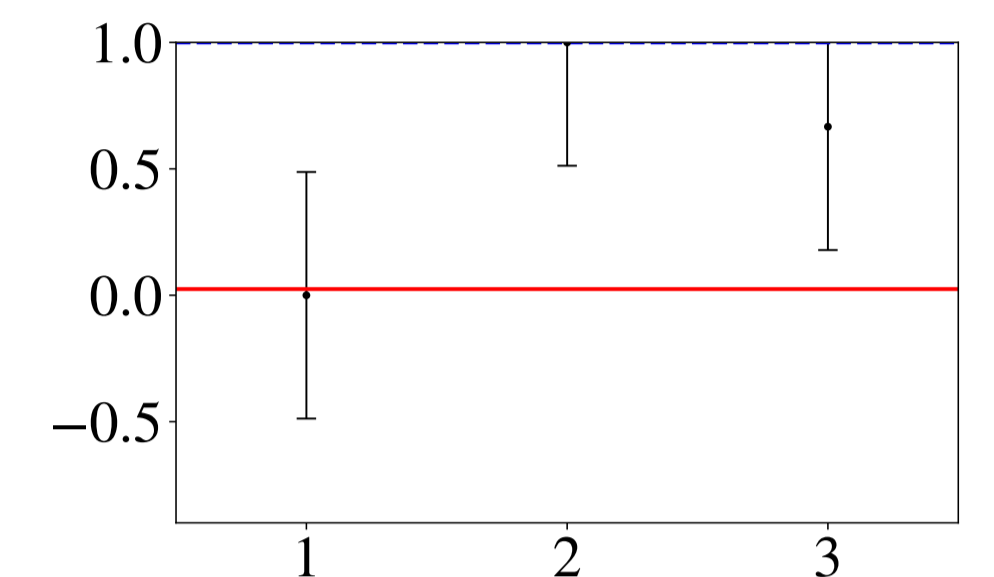
Input: K, σ, α

- 1: for $t \leftarrow K + 1, K + 2, \dots$ do
- 2: $\delta_t \leftarrow \frac{1}{Kt^\alpha}$
- 3: $h \leftarrow 1$
- 4: $\mathcal{K}_1 \leftarrow \mathcal{K}$
- 5: do
- 6: $\mathcal{K}_{h+1} \leftarrow \{i \in \mathcal{K}_h \mid \hat{\mu}_i^h(N_{i,t}) \geq \max_{j \in \mathcal{K}} \hat{\mu}_j^h(N_{j,t}) - 2c(h, \delta_t)\}$
- 7: $h \leftarrow h + 1$
- 8: while $h \leq \min_{i \in \mathcal{K}_h} N_{i,t}$
- 9: SELECT : $\{i \in \mathcal{K}_h \mid h > N_{i,t}\}$
- 10: end for



How does it work?

- ★ **Method:** Filter the set of arms
- ★ **Based on:** Expanding size of arms history (newest samples first)
- ★ **Statistical tool:** New way of using Hoeffding bound to both
 - ★ Select relevant data history
 - ★ Select arm maximizing exploration/exploitation tradeoff



Lemma

w.p. $1 - \delta_t, \forall h \leq N_{i,t}, \bar{\mu}_i^h(N_{i,t}) \geq \max_{j \in \mathcal{K}} \mu_j(N_{j,t}) - 4c(h, \delta_t)$

Guarantees

WORST-CASE BOUND

$$\mathbb{E}[R_T(\pi_F)] \leq C\sigma\sqrt{KT \log(KT)} + KL$$

PROBLEM-DEPENDENT BOUND

$$\mathbb{E}[R_T(\pi_F)] \leq \sum_{i \in \mathcal{K}} O\left(\frac{\log(KT)}{\Delta_{i,h_{i,T}^*+1}}\right)$$

$\Delta_{i,h}$ Difference between the worst reward pulled by the optimal policy and the average of the h first overpulls of arm i .

$h_{i,T}^+$ High-probability upper bound on the number of overpulls for FEWA.

Computational complexity

FEWA has an $O(t)$ time and space complexity:

1. Perform $\log(t)$ filters: $h \leftarrow 2h$
2. Keep $\log(t)$ statistics:

