



Max K-armed Bandits: on the ExtremeHunter Algorithm and an Alternative Approach

Mastane Achab¹ Stéphan Clémenton¹ Aurélien Garivier² Anne Sabourin¹ Claire Vernade¹

¹ LTCI, Télécom ParisTech, Université Paris-Saclay ² IMT, Université de Toulouse

MAX K-ARMED BANDITS

- The max K-armed bandit problem (Cicirello and Smith, 2005) is a sequential decision-making problem in an uncertain environment. At each time $t = 1, \dots, n$

- choose arm $k_t \in \{1, \dots, K\}$
- observe reward $X_{k_t, t} \sim \nu_{k_t}$.

- Objective: maximize $\mathbb{E}[\max_{1 \leq t \leq n} X_{k_t, t}]$.

- Or equivalently: minimize the *expected extreme regret*

$$\mathbb{E}[R_n] \triangleq \mathbb{E}\left[\max_{1 \leq t \leq n} X_{k^*, t}\right] - \mathbb{E}\left[\max_{1 \leq t \leq n} X_{k_t, t}\right],$$

where $k^* \triangleq \arg \max_{1 \leq k \leq K} \mathbb{E}[\max_{1 \leq t \leq n} X_{k, t}]$ is the *optimal arm*.

SECOND-ORDER PARETO

- An (α, β, C, C') -second order Pareto with cdf F verifies $\forall x \geq 0$

$$|1 - Cx^{-\alpha} - F(x)| \leq C'x^{-\alpha(1+\beta)}.$$

- As in [1], rewards $X_{k, t} \sim \nu_k$ with ν_k an $(\alpha_k, \beta_k, C_k, C')$ -second order Pareto distribution, $\alpha_k > 1$, $\beta_k > 0$, $C_k > 0$ and $C' > 0$.

Theorem 1. If $\alpha > 1$ then

$$\left| \mathbb{E}\left[\max_{1 \leq t \leq n} X_t\right] - (nC)^{1/\alpha} \Gamma(1 - 1/\alpha) \right| = \mathcal{O}\left(n^{-(\min(1, \beta) - 1/\alpha)}\right),$$

sharper than $\mathcal{O}\left(n^{\frac{1}{(1+\beta)\alpha}}\right)$ in [1].

EXTREMEETC ALGORITHM

We propose EXTREMEETC, an *Explore-Then-Commit* version of EXTREMEHUNTER [1]. Both use the *optimism-in-the-face-of-uncertainty* principle through indices

$$B_k \triangleq (n(\widehat{C}_k + \Lambda_2))^{1/\alpha_k + \Lambda_1} \Gamma(1 - 1/\alpha_k - \Lambda_1) \left(\geq \mathbb{E}\left[\max_{1 \leq t \leq T} X_{k, t}\right] \text{ with high probability} \right). \quad (1)$$

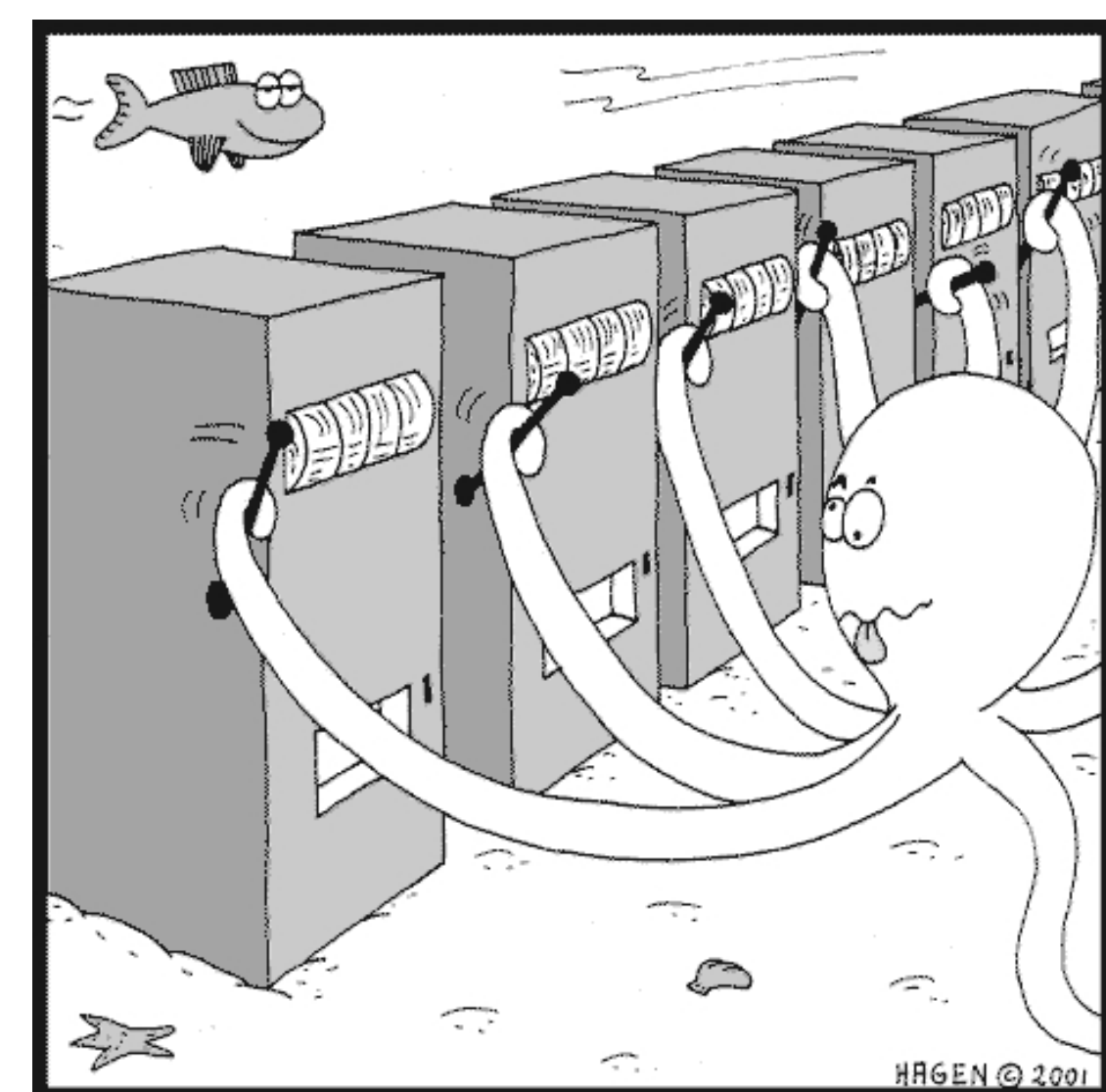
EXTREMEETC vs EXTREMEHUNTER

- Input:** n : time horizon, K : number of arms, $b > 0$ such that $b \leq \min_k \beta_k$.
- Initialize:** Pull N times each arm k and compute index B_k (see Eq. (1)).

- $k_0 = \arg \max_k B_k$
- for** $t > KN$ **do**
- for** $t > KN$ **do** Pull $k_t = \arg \max_k B_k$.
- Pull arm k_0 .
- Update index B_{k_t} .
- end for**
- end for**

Complexity Ex.ETC Ex.HUNT.

Complexity	Ex.ETC	Ex.HUNT.
Time	$\mathcal{O}(K(\log n)^{\frac{2(2b+1)}{b}})$	$\mathcal{O}(n^2)$
Memory	$\mathcal{O}(K(\log n)^{\frac{2(2b+1)}{b}})$	$\mathcal{O}(n)$



TIGHT REGRET BOUNDS

Theorem 2. (i) Upper bound for EXTREMEETC and EXTREMEHUNTER

$$\mathbb{E}[R_n] = \mathcal{O}\left((\log n)^{\frac{2(2b+1)}{b}} n^{-(1-1/\alpha_{k^*})} + n^{-(b-1/\alpha_{k^*})}\right),$$

sharper than $\mathcal{O}\left(n^{\frac{1}{(1+b)\alpha_{k^*}}}\right)$ in [1].

(ii) Lower bound for any algorithm pulling each arm at least N times

$$\mathbb{E}[R_n] = \Omega\left((\log n)^{\frac{2(2b+1)}{b}} n^{-(1-1/\alpha_{k^*})}\right).$$

When $b \geq 1$, (i) and (ii) are tight !

REDUCTION TO MULTI-ARMED BANDITS (MAB)

- Truncated rewards: $Y_{k, t} \triangleq X_{k, t} \mathbb{1}_{\{X_{k, t} > u\}}$.

$$\mathbb{E}[Y_{k, 1}] \sim_{u \rightarrow \infty} C_k \left(1 + \frac{1}{\alpha_k - 1}\right) u^{-\alpha_k + 1}.$$

- For $u > \max\left(1, \left(\frac{2C'}{C(1)}\right)^{\frac{1}{b}}, \left(\frac{3C(K)}{C(1)}\right)^{\frac{1}{\alpha(2) - \alpha(1)}}\right)$ and n large enough

$$\arg \max_{1 \leq k \leq K} \mathbb{E}[Y_{k, 1}] = \arg \min_{1 \leq k \leq K} \alpha_k = k^*.$$

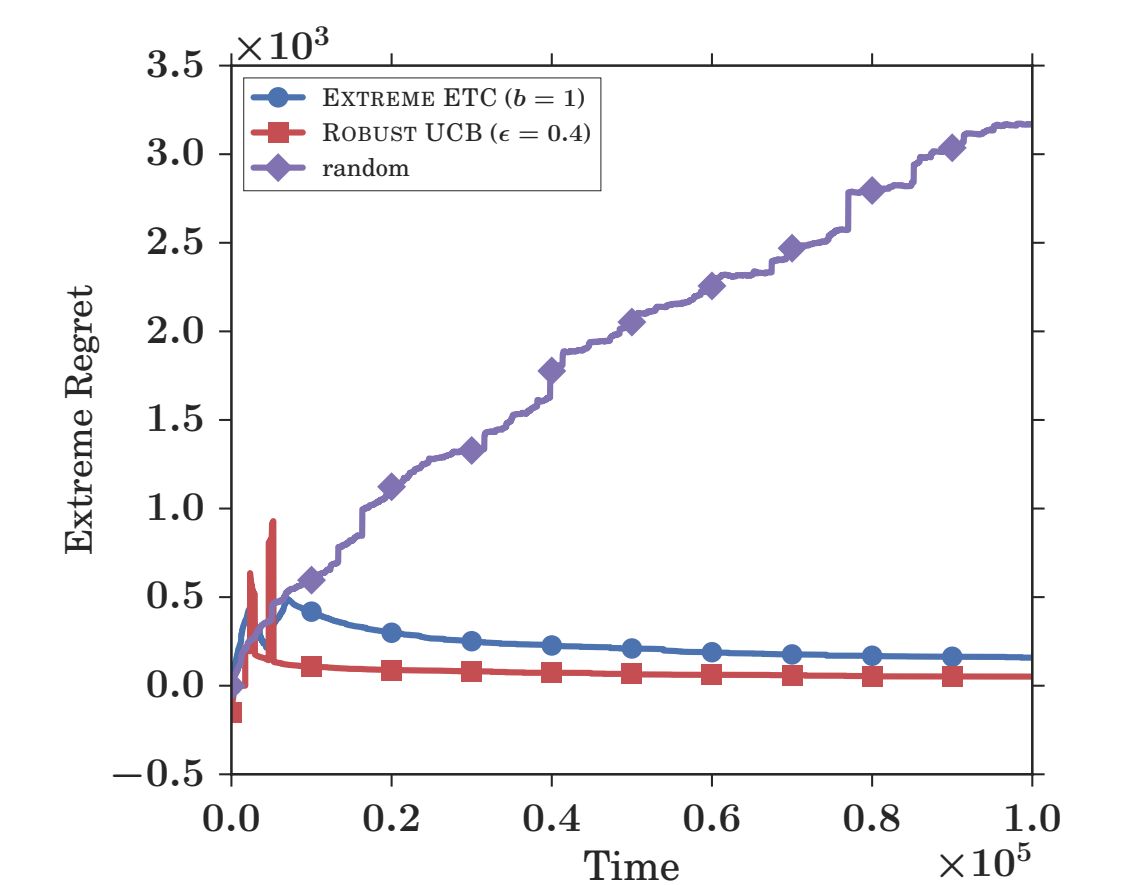
- MAB objective: maximize $\mathbb{E}[\sum_{t=1}^n Y_{k_t, t}]$.
- We use ROBUST UCB with truncated mean estimator [4]

- parameters: $\epsilon < \min_{1 \leq k \leq K} \alpha_k - 1$, $v \geq \max_{k \in [K]} \mathbb{E}[|Y_{k, 1}|^{1+\epsilon}]$
- $\mathbb{E}[\# \text{ pulls arm } k \neq k^*] = \mathcal{O}(\log n) < N$.

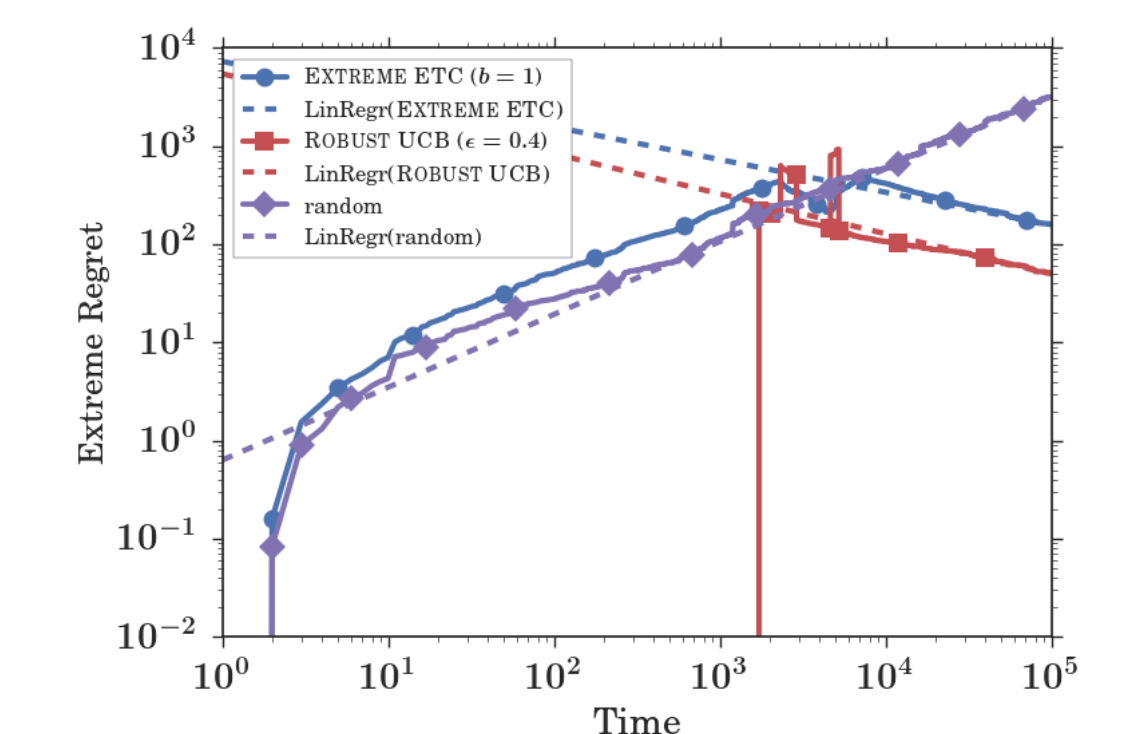
EXPERIMENTS

- time horizon $n = 10^5$
- $K = 3$ exact Pareto distributions ($\beta = +\infty$)

Arm	1	$k^* = 2$	3
α	15	1.5	10
C	10^8	1	10^5
$\mathbb{E}[X]$	3.7	3	3.5
$\mathbb{E}[\max_{1 \leq t \leq n} X_t]$	7.7	$5.8 \cdot 10^3$	11



(a) ExtremeETC vs Robust UCB vs uniform random



(b) same in log-log plot

Figure 1: Extreme regret averaged over 1000 independent trajectories.

- Fig. 1b: linear regression for EXTREMEETC over $t = 5 \cdot 10^4, \dots, 10^5$ has slope ≈ -0.333 → **validation of Theorem 2 !**

ESTIMATION OF $1/\alpha$ AND C (SEE RESP. [2] AND [3])

Assume $T \geq N \triangleq A(\log n)^{\frac{2(2b+1)}{b}}$ with b known s.t. $b \leq \beta$.

$$\widehat{1/\alpha} \triangleq \min\left(1, \left[\log\left(\frac{\sum_{t=1}^T \mathbb{1}_{\{X_t > e^r\}}}{\sum_{t=1}^T \mathbb{1}_{\{X_t > e^{r+1}\}}}\right)\right]^{-1}\right) \quad \widehat{C} \triangleq T^{-\frac{2b}{2b+1}} \sum_{t=1}^T \mathbb{1}_{\{X_t \geq T^{1/\alpha}/(2b+1)\}}$$

With probability larger than $1 - \delta$,

$$\left|\widehat{1/\alpha} - 1/\alpha\right| \leq \Lambda_1(T) \triangleq D\sqrt{\log(1/\delta)}T^{-\frac{b}{2b+1}} \quad \left|\widehat{C} - C\right| \leq \Lambda_2(T) \triangleq E\sqrt{\log(T/\delta)}\log(T)T^{-\frac{b}{2b+1}}.$$

REFERENCES

References

- Alexandra Carpentier and Michal Valko. Extreme bandits (2014).
- Alexandra Carpentier and Arlene KH Kim. Adaptive and minimax optimal estimation of the tail coefficient (2014).
- Alexandra Carpentier, Arlene KH Kim, et al. Honest and adaptive confidence interval for the tail coefficient in the pareto model (2014).
- Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail (2013).