



Max K-armed bandit: On the ExtremeHunter algorithm and beyond

Mastane Achab, Stéphan Cléménçon,
Aurélien Garivier, Anne Sabourin, Claire
Vernade

ECML PKDD 2017, Skopje





Outline

Introduction

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

EXTREMEETC algorithm

Reduction to Multi-Armed Bandits

Experiments



Outline

Introduction

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

EXTREMEETC algorithm

Reduction to Multi-Armed Bandits

Experiments

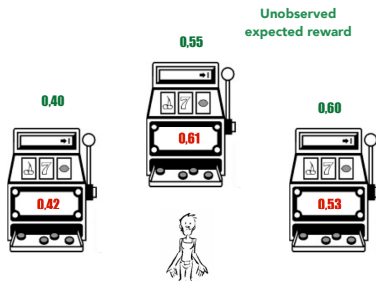
The classical Multi-Armed Bandit problem

At each time $t = 1, \dots, n$

- ▶ Play a slot machine ("pull an arm")
- ▶ Receive reward

Goal: maximize cumulative reward!

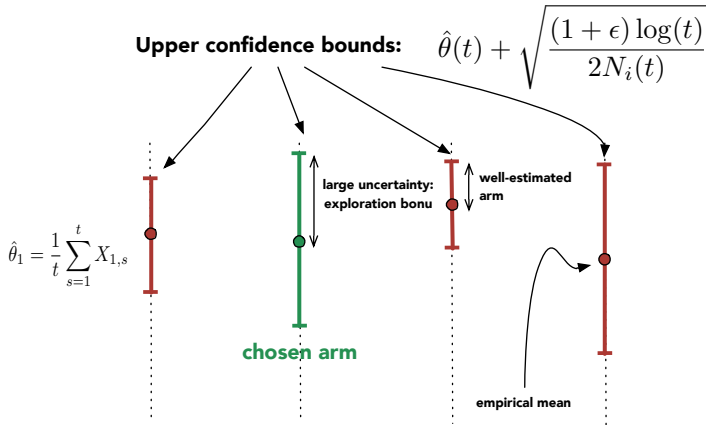
Dilemma: exploration vs exploitation



Estimated empirical
averages after a
few pulls

A successful approach: UCB algorithm (Auer et al., 2002)

- ▶ *Initialization:* pull each arm once
- ▶ Then:



The Max K-Armed Bandit problem

At each time $t = 1, \dots, n$

- ▶ Choose arm $k_t \in \{1, \dots, K\}$
- ▶ Observe reward $X_{k_t, t}$

Multi-Armed Bandits

maximize $\mathbb{E}[\sum_{t=1}^n X_{k_t, t}]$

Max K-Armed Bandits (Cicirello and Smith, 2005)

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

Extreme Regret

- ▶ optimal arm

$$k^* = \arg \max_{1 \leq k \leq K} \mathbb{E} \left[\max_{1 \leq t \leq n} X_{k,t} \right]$$

- ▶ equivalent objective

Expected extreme regret

$$\text{minimize } \mathbb{E} [R_n^\pi] = \mathbb{E} [\max_{1 \leq t \leq n} X_{k^*,t}] - \mathbb{E} [\max_{1 \leq t \leq n} X_{k_t,t}]$$

Definition

F is a 2nd-order Pareto distribution if $\forall x \geq 0$

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

with constants $\alpha, \beta, C, C' > 0$.

Some properties

- ▶ for $\beta = +\infty$, $F(x) = 1 - Cx^{-\alpha}$ (exact Pareto)
- ▶ finite moments of orders $r < \alpha$

Assumption: $\alpha > 1$ (finite mean).



Outline

Introduction

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

EXTREMEETC algorithm

Reduction to Multi-Armed Bandits

Experiments

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

$X_{1:n} \sim^{\text{iid}} 2^{\text{nd}}$ -order Pareto ($\alpha > 1, \beta, C, C'$)

$\mathbb{E}[\max_{1 \leq t \leq n} X_t] \sim_{n \rightarrow \infty} (nC)^{1/\alpha} \Gamma(1 - 1/\alpha)$ (mean of a Fréchet distribution)

Theorem 1

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

sharper than $\mathcal{O} \left(n^{\frac{1}{(1+\beta)\alpha}} \right)$ in C&V14.



Outline

Introduction

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

EXTREMEETC algorithm

Reduction to Multi-Armed Bandits

Experiments

ExtremeETC algorithm

- ▶ UCB idea:

$$\begin{aligned} B_k &= (n(\widehat{C}_k + \Lambda_2))^{\widehat{1/\alpha}_k + \Lambda_1} \Gamma(1 - \widehat{1/\alpha}_k - \Lambda_1) \\ &\geq \mathbb{E} \left[\max_{1 \leq t \leq T} X_{k,t} \right] \text{ with high probability} \end{aligned} \tag{1}$$

- ▶ Initialization: pull each arm $N = A(\log n)^{\frac{2(2b+1)}{b}}$ times

EXTREMEETC vs EXTREMEHUNTER

- 1: **Input:** n : time horizon, K : number of arms, $b > 0$ such that $b \leq \min_k \beta_k$.
- 2: **Initialize:** Pull N times each arm k and compute index B_k (see Eq. (1)).
- 3: $k_0 = \arg \max_k B_k$
- 4: **for** $t > KN$ **do**
- 5: Pull arm k_0 .
- 6: **end for**
- 3: **for** $t > KN$ **do**
- 4: Pull $k_t = \arg \max_k B_k$.
- 5: Update index B_{k_t} .
- 6: **end for**

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

Tight regret bounds

Theorem

(i) *Upper bound for*

EXTREMEETC and EXTREMEHUNTER

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

sharper than $\mathcal{O}\left(n^{\frac{1}{(1+b)\alpha_{k^*}}}\right)$ *in* C&V(14).

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

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

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

Regret bounds - idea of proof

- ▶ favorable event (\mathcal{A}): $1/\alpha_k, C_k \in$ confidence intervals
- ▶ Lemma: under (\mathcal{A}), k^* always pulled
- ▶ use Theorem 1 to control $\mathbb{E}[\max \dots]$



Outline

Introduction

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

EXTREMEETC algorithm

Reduction to Multi-Armed Bandits

Experiments

Reduction to Multi-Armed Bandits

- ▶ Idea: peak over threshold
- ▶ Truncated rewards: $Y_{k,t} = X_{k,t} \mathbb{1}_{\{X_{k,t} > u\}}$.
- ▶ $\mathbb{E}[Y_{k,1}] \sim_{u \rightarrow \infty} C_k \left(1 + \frac{1}{\alpha_k - 1}\right) u^{-\alpha_k + 1}$.
- ▶ For u and n large

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

- ▶ MAB objective: maximize $\mathbb{E} [\sum_{t=1}^n Y_{k_t,t}]$.
- ▶ We use ROBUST UCB (Bubeck et al., 2013)
 - ▶ parameters: $\epsilon < \min_{1 \leq k \leq K} \alpha_k - 1$, $v \geq \max_{k \in [K]} \mathbb{E} [|Y_{k,1}|^{1+\epsilon}]$
 - ▶ $\mathbb{E} [\# \text{ pulls arm } k \neq k^*] = \mathcal{O}(\log n) < N$.



Outline

Introduction

Controlling $\mathbb{E}[\max_{1 \leq t \leq n} X_t]$

EXTREMEETC algorithm

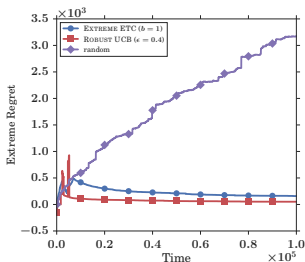
Reduction to Multi-Armed Bandits

Experiments

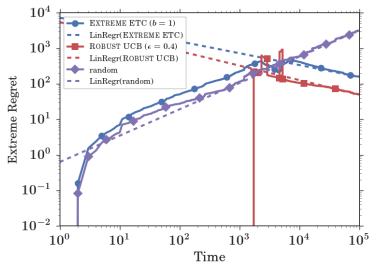
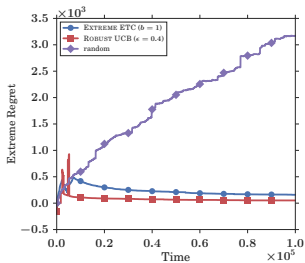
Experiments

- ▶ 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



Experiments



Linear regression for EXTREME ETC over $t = 5 \cdot 10^4, \dots, 10^5$ has slope ≈ -0.333 (with $R^2 \approx 0.97$)

→ **validation of bounds !**