



uniINF: Best-of-Both-Worlds Algorithm for Parameter-Free Heavy-Tailed MABs

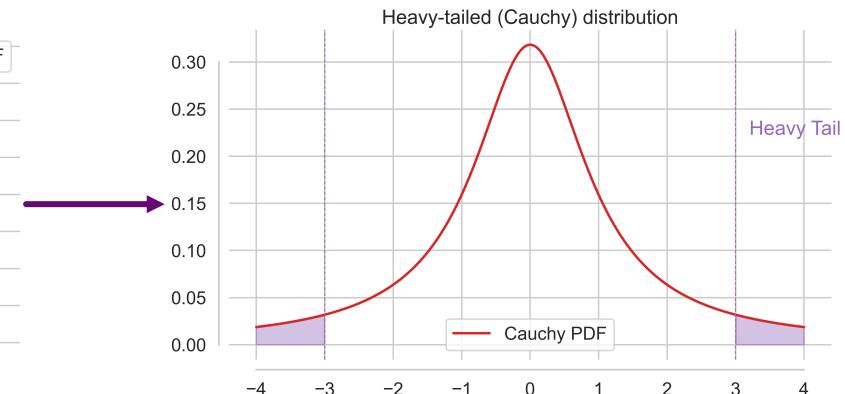
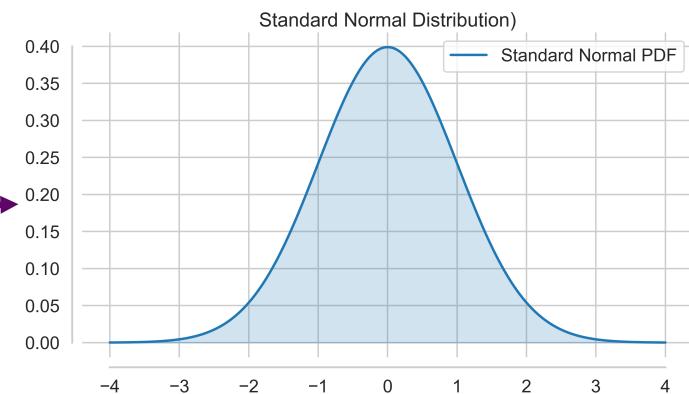
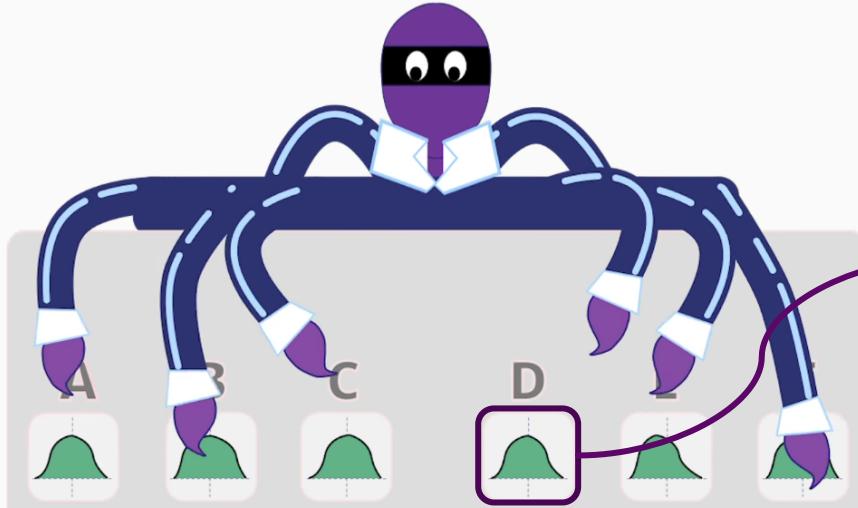
Published on ICLR2025 (Spotlight)

Yu Chen* Jiatai Huang* Yan Dai* Longbo Huang

Introduction



- Multi-Armed Bandits (MABs) is a classic theoretical model in online learning to address the **exploration-exploitation** trade-off problem.
- Many studies for MABs assume that the loss for each arm is **sub-Gaussian**, or even **bounded**.
- Challenges in the Real World: The loss distributions in many real-world scenarios are **heavy-tailed**.





- Heavy-tailed distribution ν
 - Variance of loss $l \sim \nu$ may be infinite.
 - Only α -th moment is bounded ($\alpha \in (1, 2]$): $\mathbb{E}[|l|^\alpha] \leq \sigma^\alpha$ for some $\sigma \geq 0$.
- E.g., Pareto distribution, Student t distribution.
- Empirical evidence shows that heavy-tailed distributions frequently occur in many key real-world tasks.
 - Network Routing, e.g., [\[Liebeherr et al., 2012\]](#)
 - Algorithm Portfolio Selection, e.g., [\[Gagliolo & Schmidhuber, 2011\]](#)
 - Online Deep Learning, e.g., [\[Zhang et al., 2020\]](#)



Heavy-Tailed Multi-Armed Bandits (HTMABs)

- HTMAB problem is consisted with K arms and T rounds.
- In Round $t \in [T]$,
 1. Loss vector $\mathbf{l}_t \in \mathbb{R}^K$ is generated by the environment from heavy tailed distribution ν_t where $\mathbb{E}_{l \sim \nu_{t,i}}[|l|^\alpha] \leq \sigma^\alpha$ for each arm $i \in [K]$, which is hidden to the player.
 2. Player choose an action $i_t \in [K]$ then observes and suffers a loss of l_{t,i_t} .
- The problem objective is **minimizing the regret**:

$$\mathcal{R}_T = \max_{i \in [K]} \mathbb{E} \left[\sum_{t=1}^T (l_{t,i_t} - l_{t,i}) \right]$$



Heavy-Tailed Multi-Armed Bandits (HTMABs)

- **Unknown Parameters** (σ, α) : In real applications, we usually cannot predict the heavy-tailed parameters (σ, α) of the environment, which is called “parameter free” setting.
- **Unknown Environments:** The environment may be one of two types, and the algorithm cannot predict it:
 - **Stochastic**: The loss distribution is fixed and does not change over time ($\nu_t = \nu_1$).
 - **Adversarial**: The loss can be arbitrarily generated by an adversary in advance and changes over time.
- We aim to develop the Best of Both Worlds (BoBW) algorithm for parameter-free HTMABs

Overview of Related Works For HTMABs

Our Main Contribution

- **uniINF** : the **first** Parameter-Free HTMAB algorithm that achieves the **BoBW** guarantees.
- Parameter-Free: without reliance on a prior knowledge of (σ, α) .
- BoBW: automatically achieves nearly-optimal regret bounds in both stochastic and adversarial cases.

Algorithm ^a	α -Free?	σ -Free?	Env.	Regret	Opt?
Lower Bound (Bubeck et al., 2013)	—	—	—	$\Omega\left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i}\right)^{\frac{1}{\alpha-1}} \log T\right)$ $\Omega(\sigma K^{1-1/\alpha} T^{1/\alpha})$	—
RobustUCB (Bubeck et al., 2013)	✗	✗	Only Stoc.	$\mathcal{O}\left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i}\right)^{\frac{1}{\alpha-1}} \log T\right)$ $\mathcal{O}(\sigma K^{1-1/\alpha} T^{1/\alpha})$	✓ ✓
Robust MOSS (Wei & Srivastava, 2020)	✗	✗	Only Stoc.	$\mathcal{O}\left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i}\right)^{\frac{1}{\alpha-1}} \log\left(\frac{T}{K} \left(\frac{\sigma^\alpha}{\Delta_i}\right)^{\frac{1}{\alpha-1}}\right)\right)$ $\mathcal{O}(\sigma K^{1-1/\alpha} T^{1/\alpha})$	✓ ^b ✓
APE ² (Lee et al., 2020b)	✗	✓	Only Stoc.	$\mathcal{O}\left(e^\sigma + \sum_{i \neq i^*} \left(\frac{1}{\Delta_i}\right)^{\frac{1}{\alpha-1}} (T \Delta_i^{\frac{\alpha}{\alpha-1}} \log K)^{\frac{\alpha}{(\alpha-1) \log K}}\right)$ $\tilde{\mathcal{O}}(\exp(\sigma^{1/\alpha}) K^{1-1/\alpha} T^{1/\alpha})$	✗ ✗
HTINF (Huang et al., 2022)	✗	✗	Stoc. Adv.	$\mathcal{O}\left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i}\right)^{\frac{1}{\alpha-1}} \log T\right)$ $\mathcal{O}(\sigma K^{1-1/\alpha} T^{1/\alpha})$	✓ ✓
OptHTINF (Huang et al., 2022)	✓	✓	Stoc. Adv.	$\mathcal{O}\left(\sum_{i \neq i^*} \left(\frac{\sigma^{2\alpha}}{\Delta_i^{3-\alpha}}\right)^{\frac{1}{\alpha-1}} \log T\right)$ $\mathcal{O}(\sigma^\alpha K^{\frac{\alpha-1}{2}} T^{\frac{3-\alpha}{2}})$	✗ ✗
AdaTINF (Huang et al., 2022)	✓	✓	Only Adv.	$\mathcal{O}(\sigma K^{1-1/\alpha} T^{1/\alpha})$	✓
AdaR-UCB (Genalti et al., 2024)	✓	✓	Only Stoc.	$\mathcal{O}\left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i}\right)^{\frac{1}{\alpha-1}} \log T\right)$ $\tilde{\mathcal{O}}(\sigma K^{1-1/\alpha} T^{1/\alpha})$	✓ ✓
uniINF (Ours)	✓	✓	Stoc. Adv.	$\mathcal{O}\left(K \left(\frac{\sigma^\alpha}{\Delta_{\min}}\right)^{\frac{1}{\alpha-1}} \log T \cdot \log \frac{\sigma^\alpha}{\Delta_{\min}}\right)$ $\tilde{\mathcal{O}}(\sigma K^{1-1/\alpha} T^{1/\alpha})$	✓ ^c ✓

^a α -Free? and σ -Free? denotes whether the algorithm is parameter-free w.r.t. α and σ , respectively. Env. includes the environments that the algorithm can work; if one algorithm can work in **both** stochastic and adversarial environments, then we mark this column by green. Regret describes the algorithmic guarantees, usually (if applicable) instance-dependent ones above instance-independent ones. Opt? means whether the algorithm matches the instance-dependent lower bound by Bubeck et al. (2013) up to constant factors, or the instance-independent lower bound up to logarithmic factors.

^bUp to $\log(\sigma^\alpha)$ and $\log(1/\Delta_i^\alpha)$ factors.

^cUp to $\log(\sigma^\alpha)$ and $\log(1/\Delta_{\min})$ factors when all Δ_i 's are similar to the dominant sub-optimal gap Δ_{\min} .



Main Results

■ Assumption 0 [Unique Best Arm, very common in MAB literatures]:

There exists an optimal arm $i^* \in [K]$ such that $\mathbb{E}_{l \sim \nu_{t,i}}[l] - \mathbb{E}_{l \sim \nu_{t,i^*}}[l] > 0$ for any $i \neq i^*$.

■ Assumption I [Truncated Non-Negativity, Huang et al., 2022]:

For any t and optimal loss $l \sim \nu_{t,i^*}$, $\mathbb{E}[l \cdot \mathbf{1}[|l| \geq M]] \geq 0$ holds for any $M \geq 0$.

With Ass 0 and I, we have the following main theorem: **uniINF is BoBW algorithm.**

Main Theorem: uniINF guarantees near-optimal regret automatically under **both** stochastic and adversarial environments **without any prior knowledge** of the heavy-tailed parameters α and σ .

- ▶ **Adversarial** environments: uniINF achieves $\mathcal{R}_T = \tilde{\mathcal{O}}\left(\sigma K^{1-1/\alpha} T^{1/\alpha}\right)$.
 - ▶ Matches stochastic-case instance-independent lower bound (Bubeck et al., 2013) up to logs.
- ▶ **Stochastic** environments: uniINF achieves $\mathcal{R}_T = \tilde{\mathcal{O}}\left(K\left(\frac{\sigma^\alpha}{\Delta_{\min}}\right)^{\frac{1}{\alpha-1}} \log T \cdot \log \frac{\sigma^\alpha}{\Delta_{\min}}\right)$.
 - ▶ Nearly matches stochastic-case instance-dependent lower bound up to polylog($\sigma^\alpha, \Delta_{\min}^{-1}$).



Discussion on the Assumption I

- We summarize the existing results in HTMABs for the relationship between the prior knowledge of heavy-tailed parameters (σ, α) and Assumption 1 in the following table.
 - [Cheng et al. (2024)]: when parameters (σ, α) are known, Assumption 1 is redundant.
 - [Genalti et al. (2024)]: achieving optimal regret guarantees is impossible without any assumptions.

Assumptions	Known (σ, α)	Unknown (σ, α)
With Assumption 1	HTINF achieves BoBW (Huang et al., 2022)	uni INF achieves BoBW (Our Work)
Weaker than Assumption 1?	Inferred by Cheng et al. (2024) (discussed below)	Open Problem
Without Any Assumption	SAO-HT achieves BoBW (Cheng et al., 2024)	No BoBW possible (Genalti et al., 2024, Theorems 2 & 3)

- **Our work:** Assumption 1 is sufficient for achieving BoBW under parameter-free settings.

uniINF: The Universal INF-type Algorithm For Parameter-Free HTMAB



■ **Core framework:** "Follow the Regularized Leader" (FTRL).

■ Three key-innovations components:

I. Refined Analysis for Log-Barrier Regularizer

2. Adaptive Skipping-Clipping Loss Tuning

3. Auto-Balancing Learning Rates

Algorithm uniINF: the universal INF-type algorithm for Parameter-Free HTMAB

- 1: Initialize the learning rate $S_1 \leftarrow 4$.
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: Apply Follow-the-Regularized-Leader (FTRL) to calculate the action $\mathbf{x}_t \in \Delta^{[K]}$ with the log-barrier regularizer Ψ_t defined in Equation (1): \triangleright FTRL with refined log-barrier regularizer.

$$\mathbf{x}_t \leftarrow \operatorname{argmin}_{\mathbf{x} \in \Delta^{[K]}} \left(\sum_{s=1}^{t-1} \langle \tilde{\ell}_s, \mathbf{x} \rangle + \Psi_t(\mathbf{x}) \right), \quad \Psi_t(\mathbf{x}) := -S_t \sum_{i=1}^K \log x_i \quad (1)$$

- 4: Sample action $i_t \sim \mathbf{x}_t$. Play i_t and observe feedback ℓ_{t,i_t} .
- 5: **for** $i = 1, 2, \dots, K$ **do**
- 6: Calculate $\ell_{t,i}^{\text{skip}}$ and $\ell_{t,i}^{\text{clip}}$ with action-dependent threshold $C_{t,i}$ by Equation (2). \triangleright Adaptive skipping-clipping loss tuning.

$$\ell_{t,i}^{\text{skip}} := \begin{cases} \ell_{t,i} & \text{if } |\ell_{t,i}| < C_{t,i} \\ 0 & \text{otherwise} \end{cases}, \quad \ell_{t,i}^{\text{clip}} := \begin{cases} C_{t,i} & \text{if } \ell_{t,i} \geq C_{t,i} \\ -C_{t,i} & \text{if } \ell_{t,i} \leq -C_{t,i} \\ \ell_{t,i} & \text{otherwise} \end{cases}, \quad \text{with } C_{t,i} := \frac{S_t}{4(1-x_{t,i})} \quad (2)$$

- 7: Calculate the importance sampling estimate of ℓ_t^{skip} , namely $\tilde{\ell}_t$, where

$$\tilde{\ell}_t = \frac{\ell_{t,i_t}^{\text{skip}}}{x_{t,i_t}} \cdot \mathbb{1}[i = i_t], \quad \forall i \in [K].$$

- 8: Update the learning rate S_{t+1} as \triangleright Auto-balancing learning rates

$$S_{t+1}^2 = S_t^2 + (\ell_{t,i_t}^{\text{clip}})^2 \cdot (1 - x_{t,i_t})^2 \cdot (K \log T)^{-1}. \quad (3)$$



Analysis Preparation: FTRL Regret Decomposition

- By basic FTRL regret decomposition, we have

$$\mathcal{R}_T \leq \underbrace{\mathbb{E} \left[\sum_{t=1}^T \text{DIV}_t \right]}_{\text{BREGMAN DIVERGENCE TERMS}} + \underbrace{\mathbb{E} \left[\sum_{t=0}^{T-1} \text{SHIFT}_t \right]}_{\Psi\text{-SHIFTING TERMS}} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T |\text{SKIPERR}_t| \cdot \mathbb{1}[i_t \neq i^*] \right]}_{\text{SUB-OPTIMAL SKIPPING LOSSES}} + \sigma K.$$

1. Bregman Divergence terms

$$\text{Div}_t = \Psi_t(\mathbf{x}_t) - \Psi_t(\mathbf{z}_t) - \langle \nabla \Psi_t(\mathbf{x}_t), \mathbf{z}_t - \mathbf{x}_t \rangle,$$

where $\mathbf{z}_t := \operatorname{argmin}_{\mathbf{z} \in \Delta^{[K]}} (\sum_{s=1}^t \langle \tilde{\mathbf{l}}_s, \mathbf{z} \rangle + \Psi_t(\mathbf{z}))$,

2. The Ψ -Drift term Shift_t for benchmark $\tilde{\mathbf{y}}$, ($\tilde{y}_{i^*} = \frac{1}{T}$, $\tilde{y}_i = 1 - \frac{K-1}{T}$, $\forall i \neq i^*$)

$$\text{Shift}_t = (\Psi_{t+1} - \Psi_t)(\tilde{\mathbf{y}}) - (\Psi_{t+1} - \Psi_t)(\mathbf{x}_{t+1}),$$

3. Skipping Error terms SkipErr_t :

$$\text{SkipErr}_t = l_{t,i_t} - l_{t,i_t}^{\text{skip}}.$$

Key Technique: Refined Analysis for Log-Barrier Regularizer



- **Log-barrier regularizer** $\Psi_t(\mathbf{x}) = -S_t \sum_{i=1}^K \log(x_i)$, where S_t^{-1} is learning rate in round t .
- It is known that log-barrier applied to a loss sequence $\{\mathbf{l}_t\}_{t=1}^T$ ensures the control of Div_t
 - [Foster et al., 2016, Lemma 16] $\text{Div}_t \leq S_t^{-1} \sum_i x_{t,i}^2 l_{t,i}^2$ for non-negative \mathbf{l}_t .
 - [Dai et al., 2023, Lemma 3.1] $\text{Div}_t \leq S_t^{-1} \sum_i x_{t,i} l_{t,i}^2$ for general \mathbf{l}_t .
- **Our work:** $\text{Div}_t \leq S_t^{-1} \sum_i x_{t,i}^2 (1 - x_{t,i})^2 l_{t,i}^2$ for general \mathbf{l}_t when S_t is adequately large compared to $\|\mathbf{l}_t\|_\infty$.
- Extra $(1 - x_{t,i})^2$ terms : contributes to exclude the optimal arm $i^* \in [K]$ in stochastic case analysis.



Key Technique: Adaptive Skipping-Clipping Loss Tuning

- We cannot straightforward control the Bregman divergence term

$$\text{Div}_t \leq S_t^{-1} \sum_i x_{t,i}^2 (1 - x_{t,i})^2 l_{t,i}^2$$

with original loss vector \mathbf{l}_t , since $\mathbb{E}[l_{t,i}^2]$ might be unbounded in HTMABs.

- **Previous Solution:** [Huang et al., 2022] developed a skipping technique that replaces the extremely large loss with 0 with dynamic threshold.
- **Challenge:** Skipping is not a free lunch, which will influence the update of learning rates and reduce the learning efficiency.



Key Technique: Adaptive Skipping-Clipping Loss Tuning

- **Our Solution:** we apply skipping-clipping techniques with adaptive thresholds $C_{t,i}$.
- We tune loss vector ℓ_t via Eq.(2):

$$\ell_{t,i}^{\text{skip}} := \begin{cases} \ell_{t,i} & \text{if } |\ell_{t,i}| < C_{t,i} \\ 0 & \text{otherwise} \end{cases}, \quad \ell_{t,i}^{\text{clip}} := \begin{cases} C_{t,i} & \text{if } \ell_{t,i} \geq C_{t,i} \\ -C_{t,i} & \text{if } \ell_{t,i} \leq -C_{t,i} \\ \ell_{t,i} & \text{otherwise} \end{cases}, \quad \text{with } C_{t,i} := \frac{S_t}{4(1-x_{t,i})} \quad (2)$$

- This adaptive skipping-clipping technique ensure that **every loss will influence the learning process**.



Key Technique: Auto-Balancing Learning Rates

- **Challenge:** How to set optimal learning rates S_t when (σ, α) is unknown?
- **Our Solution:** Automatically balance Div_t and Shift_t by proper S_t

$$\underbrace{\text{Div}_t \leq \mathcal{O} \left(S_t^{-1} \left(\ell_{t,i_t}^{\text{clip}} \right)^2 (1 - x_{t,i_t})^2 \right)}_{\text{Bregman Divergence}}, \quad \underbrace{\text{Shift}_t \leq \mathcal{O} ((S_{t+1} - S_t) \cdot K \log T)}_{\Psi\text{-Shifting}}. \quad (4)$$

- Thus, to make Div_t roughly the same as Shift_t , it suffices to ensure

$$(S_{t+1} - S_t)S_t \approx \left(\ell_{t,i_t}^{\text{clip}} \right)^2 (1 - x_{t,i_t})^2 \cdot (K \cdot \log T)^{-1}$$

- Moreover, notice that $(S_{t+1} - S_t)S_t \approx S_{t+1}^2 - S_t^2$ if S_t does not change too much. Therefore, our definition of S_t follows

$$S_{t+1}^2 = S_t^2 + (\ell_{t,i_t}^{\text{clip}})^2 \cdot (1 - x_{t,i_t})^2 \cdot (K \log T)^{-1}. \quad (3)$$



- **Our Contribution:** The First Parameter-Free BoBW Algorithm for HTMABs.
- We introduced **uniINF**, the **first** algorithm that solves the Parameter-Free Heavy-Tailed Multi-Armed Bandit (HTMAB) problem with the Best-of-Both-Worlds (BoBW) property.
 - uniINF performs nearly-optimally in both **stochastic** and **adversarial** environments without requiring any prior knowledge of the environment type or its heavy-tail parameters (σ, α) .
- Our analytic tools and key innovation techniques could be independent interests in MABs community.



Thanks for Listening