

清华大学 交叉信息研究院
Institute for Interdisciplinary Information Sciences, Tsinghua University

Adaptive Best-of-Both-Worlds Algorithm for Heavy-Tailed Multi-Armed Bandits

Longbo Huang
Tsinghua University

The Multi-Armed Bandit Model

The MAB framework has many applications

- Online advertising
- Wireless communication
- Clinical trial
- Recommender system
- AI technology
- Finance
- ...

Extensively studied

[Auer et al. 2002a][Auer et al 2002b][Agrawal&Goyal 2012][Bubeck&Cesa-Bianchi 2012][Tao et al. 2018][Kuroki et al. 2020][Du et al. 2020] [Wang&Chen 2022][Chen&Zhao&Li 2022] ...

Multi-Armed Bandits

The MAB problem

- T time steps, K actions (“arms”)
- $\{l_{t,a}\}$: $T \times K$ loss matrix
- Each time we choose A_t , suffer & observe a loss l_{t,A_t}
- Minimize “pseudo-regret”

$$\max_{a \in [K]} \mathbb{E} \left[\sum_{t=1}^T l_{t,A_t} - \sum_{t=1}^T l_{t,a} \right]$$

Heavy-Tailed Multi-Armed Bandits

The Heavy-Tailed MAB problem

- T time steps, K actions (“arms”)
- $\{l_{t,a}\}$: $T \times K$ loss matrix
- Adversary picks (α, σ) -heavy-tailed distributions $\nu_{t,1}, \dots, \nu_{t,K}$ with
$$\mathbb{E}_{X \sim \nu}[|X|^\alpha] \leq \sigma^\alpha, 1 < \alpha \leq 2$$
- Each time we choose A_t , suffer & observe a loss $l_{t,A_t} \sim \nu_{t,A_t}$
- Minimize “pseudo-regret”

$$\max_{a \in [K]} \mathbb{E} \left[\sum_{t=1}^T l_{t,A_t} - \sum_{t=1}^T l_{t,a} \right]$$

Heavy-Tailed Multi-Armed Bandits

The Heavy-Tailed MAB problem

- T time steps, K actions (“arms”)
- $\{l_{t,a}\}$: $T \times K$ loss matrix
- Adversary picks (α, σ) -heavy-tailed distributions $\nu_{t,1}, \dots, \nu_{t,K}$ with

$$\mathbb{E}_{X \sim \nu}[|X|^\alpha] \leq \sigma^\alpha, 1 < \alpha \leq 2$$

Natural generalization in both directions

- Stochastic heavy-tailed MAB: $\nu_{1,a} = \nu_{2,a} = \dots = \nu_{T,a}$
- Classical Adversarial MAB: $\nu_{t,a}$ is a Dirac-measure at $l_{t,a} \in [0,1]$

Heavy-Tailed Multi-Armed Bandits

The Heavy-Tailed MAB problem

- T time steps, K actions (“arms”)
- $\{l_{t,a}\}$: $T \times K$ loss matrix
- Adversary picks (α, σ) -heavy-tailed distributions $\nu_{t,1}, \dots, \nu_{t,K}$ with

$$\mathbb{E}_{X \sim \nu}[|X|^\alpha] \leq \sigma^\alpha, 1 < \alpha \leq 2$$

Challenges

- Potentially unbounded 2nd moment (estimation and concentration)
- Unknown σ and α values (learning)

Representative Results on Heavy-Tailed MAB

Algorithm	Loss Type	Prior Knowledge	Total Regret
Lower-bounds (Bubeck et al., 2013)	Stochastic ^a	α, σ	$\Omega \left(\sigma^{\frac{\alpha}{\alpha-1}} \sum_{i \neq i^*} \Delta_i^{-\frac{1}{\alpha-1}} \log T \right)$
			$\Omega \left(\sigma K^{1-1/\alpha} T^{1/\alpha} \right)$
RobustUCB (Bubeck et al., 2013)	Stochastic	α, σ	$\mathcal{O} \left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i} \right)^{\frac{1}{\alpha-1}} \log T \right)$ (optimal)
			$\mathcal{O} \left(\sigma (K \log T)^{1-1/\alpha} T^{1/\alpha} \right)$ (sub-optimal for $\log T$ factors)
Lee et al. (2020)	Stochastic	α ; require $\mu_i \in [0, 1]$	$\mathcal{O} \left(K^{1-1/\alpha} T^{1/\alpha} \log K \right)^b$ (sub-optimal for $\log K$ factors)
$1/2$ -Tsallis-INF (Zimmert & Seldin, 2019)	SCA-unique ^c	$[0, 1]$ -bounded losses	$\mathcal{O} \left(\sum_{i \neq i^*} \frac{1}{\Delta_i} \log T \right)$ (optimal for $\alpha = 2, \sigma = 1$ case)
	Adversarial		$\mathcal{O} \left(\sqrt{KT} \right)$ (optimal for $\alpha = 2, \sigma = 1$ case)
HTINF (ours)	SCA-unique	α, σ	$\mathcal{O} \left(\sum_{i \neq i^*} \left(\frac{\sigma^\alpha}{\Delta_i} \right)^{\frac{1}{\alpha-1}} \log T \right)$ (optimal)
	Adversarial		$\mathcal{O} \left(\sigma K^{1-1/\alpha} T^{1/\alpha} \right)$ (optimal)
Optimistic HTINF (ours)	SCA-unique	None	$\mathcal{O} \left(\sum_{i \neq i^*} \left(\frac{\sigma^{2\alpha}}{\Delta_i^{3-\alpha}} \right)^{\frac{1}{\alpha-1}} \log T \right)$
	Adversarial		$\mathcal{O} \left(\sigma^\alpha K^{\frac{\alpha-1}{2}} T^{\frac{3-\alpha}{2}} \right)$
AdaTINF (ours)	Adversarial	None ^d	$\mathcal{O} \left(\sigma K^{1-1/\alpha} T^{1/\alpha} \right)$ (optimal)

Need to know σ, α before-hand

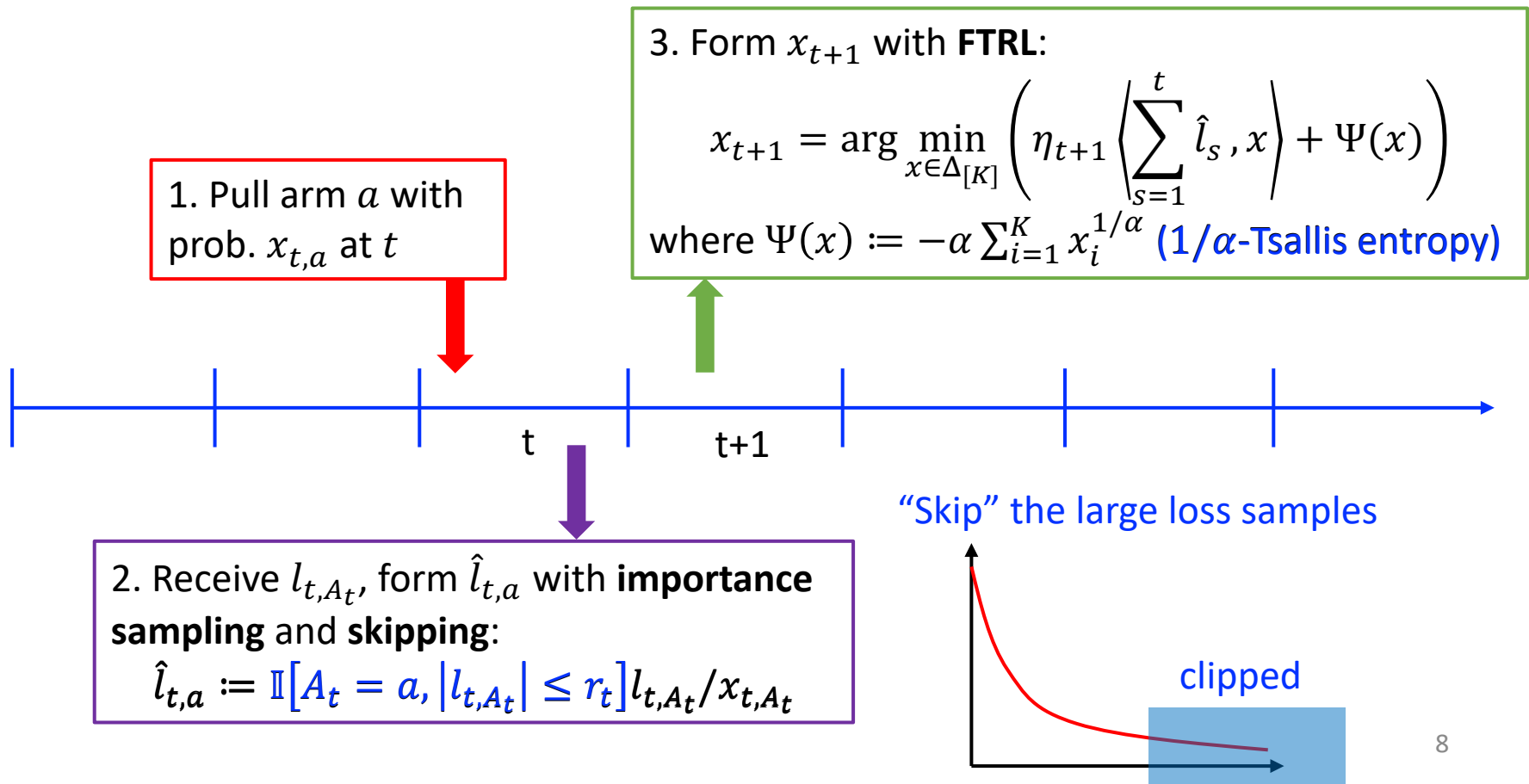
Our Contributions

Three novel algorithms

- Heavy-Tail Tsallis-INF (**HTINF**) – known σ, α
 - First to achieve *best-of-both-worlds* for heavy-tailed MAB
 - Applicable to unknown σ, α case (**OptTINF**): $O(\log T)$ for stochastic and $O\left(K^{\frac{\alpha-1}{2}} T^{\frac{3-\alpha}{2}}\right)$ for adversarial
- Adaptive Tsallis-INF (**AdaTINF**) – zero knowledge
 - Optimal $O\left(\sigma K^{1-\frac{1}{\alpha}} T^{\frac{1}{\alpha}}\right)$ regret for adversarial

Heavy-Tail Tsallis-INF (HTINF) – Known σ, α

- Based on Follow-the-Regularized-Leader (FTRL)
- A novel skipping idea to “clip” large samples



Heavy-Tail Tsallis-INF (HTINF) – Known σ, α

Algorithm 1 Heavy-Tail Tsallis-INF (HTINF)

Input: Number of arms K , heavy-tail parameters α and σ

Output: Sequence of actions $i_1, i_2, \dots, i_T \in [K]$

- 1: **for** $t = 1, 2, \dots$ **do**
 - 2: Calculate policy with learning rate $\eta_t^{-1} = \sigma t^{1/\alpha}$;
 Pick the regularizer $\Psi(x) = -\alpha \sum_{i=1}^K x_i^{1/\alpha}$ ← $1/\alpha$ -Tsallis entropy
 - $$x_t \leftarrow \operatorname{argmin}_{x \in \Delta_{[K]}} \left(\eta_t \sum_{s=1}^{t-1} \langle \hat{\ell}_s, x \rangle + \Psi(x) \right)$$
 ← Follow-the-Regularized-Leader
 - 3: Sample new action $i_t \sim x_t$.
 - 4: Calculate the skipping threshold $r_t \leftarrow \Theta_\alpha \eta_t^{-1} x_{t,i_t}^{1/\alpha}$ ← Skipping threshold (avoid overly large loss in estimation)
 where $\Theta_\alpha = \min\{1 - 2^{-\frac{\alpha-1}{2\alpha-1}}, (2 - \frac{2}{\alpha})^{\frac{1}{2-\alpha}}\}$.
 - 5: Play according to i_t and observe loss feedback ℓ_{t,i_t} .
 - 6: **if** $|\ell_{t,i_t}| > r_t$ **then**
 - 7: $\hat{\ell}_t \leftarrow \mathbf{0}$.
 - 8: **else**
 - 9: Construct weighted importance sampling loss estimator $\hat{\ell}_{t,i} \leftarrow \frac{\ell_{t,i}}{x_{t,i}} \mathbb{1}[i = i_t], \forall i \in [K]$. ← Importance Sampling (biased due to skipping)
 - 10: **end if**
 - 11: **end for**
-

Heavy-Tail Tsallis-INF (HTINF) – Known σ, α

Theorem (Informal) HTINF achieves

- Adversarial environment

$$R_T \leq O\left(K^{1-\frac{1}{\alpha}} T^{\frac{1}{\alpha}} \log K\right)$$

- Stochastic environment

$$R_T \leq O\left(\sigma^{\frac{\alpha}{\alpha-1}} \sum_{a \neq a^*} \Delta_i^{-\frac{1}{\alpha-1}} \log T\right)$$

Remark

- **Best-of-both-worlds:** Both cases are optimal without knowing which environment beforehand

Adaptivity to unknown σ, α

Theorem (Informal) When σ, α are unknown, running HTINF with $\sigma = 1, \alpha = 2$ (OptTINF) achieves

- Adversarial environment

$$R_T \leq O\left(\sigma^\alpha K^{\frac{\alpha-1}{2}} T^{\frac{3-\alpha}{2}} + \sqrt{KT}\right)$$

- Stochastic environment

$$R_T \leq O\left(\sigma^{\frac{2\alpha}{\alpha-1}} \sum_{a \neq a^*} \Delta_i^{-\frac{3-\alpha}{\alpha-1}} \log T\right)$$

Remarks

- Still $O(\log T)$ regret for stochastic case
- $o(T)$ regret for adversarial case

HTINF Regret Analysis

Regret decomposition of HTINF

$$\begin{aligned}\mathcal{R}_T(y) &\triangleq \sum_{t=1}^T \mathbb{E}[\langle x_t - y, \mu_t \rangle] \quad (y \in \Delta_{[K]}) \\ &= \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle x_t - y, \mu_t - \mu'_t \rangle \right]}_{\text{Skipping gap}} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle x_t - y, \hat{\ell}_t \rangle \right]}_{\text{FTRL Error}}\end{aligned}$$

where $\mu'_{t,i} \triangleq \mathbb{E}[\ell_{t,i} \mathbb{1}[|\ell_{t,i}| \leq r_t] \mid \mathcal{F}_{t-1}, i_t = i]$ is the clipped expectation

Skipping threshold impacts the regret

- A larger r_t leads to a smaller skipping gap but a larger FTRL error
- Optimal tradeoff achieved at $r_t = \Theta(\eta_t^{-1} x_{t,i_t}^{1/\alpha})$

HTINF Regret Analysis

Regret decomposition of HTINF

$$\begin{aligned}\mathcal{R}_T(y) &\triangleq \sum_{t=1}^T \mathbb{E}[\langle x_t - y, \mu_t \rangle] \quad (y \in \Delta_{[K]}) \\ &= \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle x_t - y, \mu_t - \mu'_t \rangle \right]}_{\text{Skipping gap}} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle x_t - y, \hat{\ell}_t \rangle \right]}_{\text{FTRL Error}}\end{aligned}$$

where $\mu'_{t,i} \triangleq \mathbb{E}[\ell_{t,i} \mathbb{1}[|\ell_{t,i}| \leq r_t] \mid \mathcal{F}_{t-1}, i_t = i]$ is the clipped expectation

Analysis idea

- “Self-bounding property” similar to (Zimmert & Seldin, 2019)
- First result for $1/\alpha$ -Tsallis ($\alpha < 2$) regularized FTRL

HTINF Regret Analysis

Regret decomposition of HTINF

$$\begin{aligned}
 \mathcal{R}_T(y) &= \mathbb{E} \left[\underbrace{\sum_{t=1}^T \langle x_t - y, \mu_t - \mu'_t \rangle}_{\text{By choice of } r_t} \right] + \mathbb{E} \left[\underbrace{\sum_{t=1}^T \langle x_t - y, \hat{\ell}_t \rangle}_{\text{FTRL Error}} \right] \\
 &\leq 10\sigma(T+1)^{1/\alpha} K^{1-1/\alpha} \quad \parallel
 \end{aligned}$$

Adversarial case

Instance-independent

$$R_T \leq O\left(\sigma^\alpha K^{\frac{\alpha-1}{2}} T^{\frac{3-\alpha}{2}}\right)$$

$$\begin{aligned}
 &\underbrace{\sum_{t=1}^T \eta_t^{-1} D_\Psi(x_t, z_t)}_{\text{By Bregman divergence}} \quad \text{Part B} \quad + \quad \underbrace{\sum_{t=1}^T (\eta_t^{-1} - \eta_{t-1}^{-1}) (\Psi(y) - \Psi(x_t))}_{\text{By Tsallis Entropy}} \quad \text{Part A} \\
 &\leq 8\sigma t^{1/\alpha-1} K^{1-1/\alpha} \quad \leq 4\sigma(T+1)^{1/\alpha} K^{1-1/\alpha}
 \end{aligned}$$

HTINF Regret Analysis

Regret decomposition of HTINF

$$\begin{aligned}
 \mathcal{R}_T(y) &= \mathbb{E} \left[\underbrace{\sum_{t=1}^T \langle x_t - y, \mu_t - \mu'_t \rangle}_{\text{Skipping gap}} \right] + \mathbb{E} \left[\underbrace{\sum_{t=1}^T \langle x_t - y, \hat{\ell}_t \rangle}_{\text{FTRL Error}} \right] \\
 &\leq 5\sigma \sum_{t=1}^T \sum_{i \neq i^*} t^{1/\alpha-1} x_{t,i}^{1/\alpha} \quad \parallel
 \end{aligned}$$

Stochastic case

Instance-dependent

$$R_T \leq O(\log T)$$

$$\begin{aligned}
 &\underbrace{\sum_{t=1}^T \eta_t^{-1} D_{\Psi}(x_t, z_t)}_{\text{Part B}} + \underbrace{\sum_{t=1}^T (\eta_t^{-1} - \eta_{t-1}^{-1}) (\Psi(y) - \Psi(x_t))}_{\text{Part A}} \\
 &\leq 8\sigma t^{1/\alpha-1} \sum_{i \neq i^*} x_{t,i}^{1/\alpha} \quad \text{By Bregman divergence} \quad \leq \sum_{t=1}^T 2\sigma t^{1/\alpha-1} \sum_{i \neq i^*} x_{t,i}^{1/\alpha} \quad \text{By Tsallis Entropy}
 \end{aligned}$$

Adaptive Tsallis-INF (AdaTINF) – Unknown σ, α

Now consider not having σ, α

$$\mathcal{R}_T(y) = \mathbb{E} \left[\sum_{t=1}^T \langle x_t - y, \mu_t - \mu'_t \rangle \right] + \sum_{t=1}^T \eta_t^{-1} D_{\Psi}(x_t, z_t) + \sum_{t=1}^T (\eta_t^{-1} - \eta_{t-1}^{-1}) (\Psi(y) - \Psi(x_t))$$

Skipping gap

Decreasing in skipping threshold

Part B

Increasing in skipping threshold

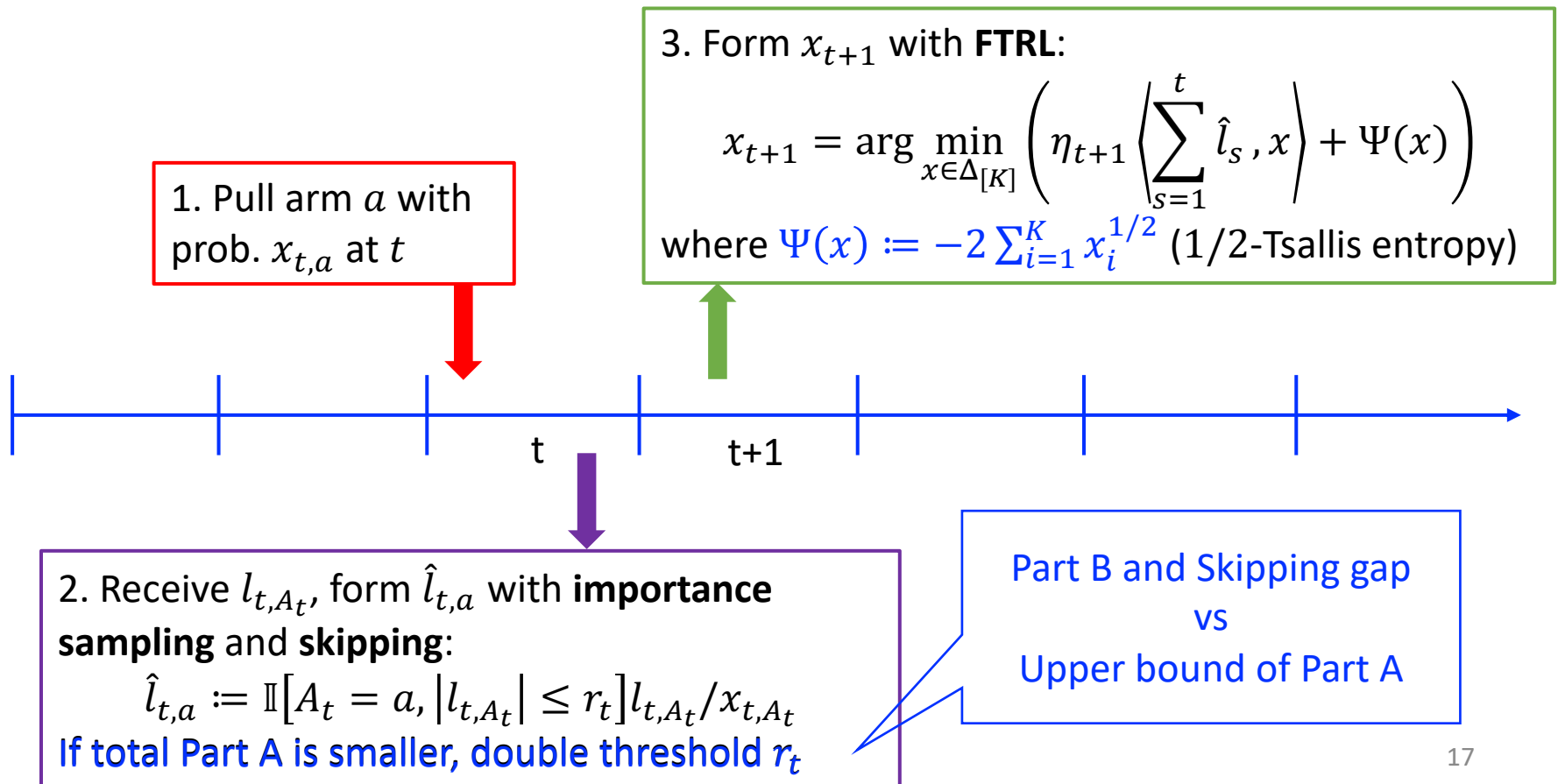
Part A

FTRL Error = Part A + Part B

AdaTINF Idea: adjust the tradeoff between Part A vs Skipping + Part B at runtime

Adaptive Tsallis-INF (AdaTINF) – Unknown σ, α

AdaTINF Idea: Using doubling trick to tune learning rate and skipping threshold



Adaptive Tsallis-INF (AdaTINF) – Unknown σ, α

Theorem (Informal) AdaTINF achieves the following for the adversarial environment:

$$R_T \leq O\left(\sigma K^{1-\frac{1}{\alpha}} T^{\frac{1}{\alpha}}\right)$$

Remarks

- **Minimax optimal:** Matches the lower bound (Bubeck et al 2013)
- Prior concentration methods heavily rely on knowing α
- Achieving instance-dependent optimality and BoBW are still open

Conclusions

Three novel algorithms for heavy-tailed MAB

- Heavy-Tail Tsallis-INF (**HTINF**) – known σ, α
 - First to achieve *best-of-both-worlds* for heavy-tailed MAB
 - Applicable to unknown σ, α case (**OptTINF**): $O(\log T)$ for stochastic and $O\left(K^{\frac{\alpha-1}{2}} T^{\frac{3-\alpha}{2}}\right)$ for adversarial
- Adaptive Tsallis-INF (**AdaTINF**) – zero knowledge
 - Optimal $O\left(\sigma K^{1-\frac{1}{\alpha}} T^{\frac{1}{\alpha}}\right)$ regret for adversarial

Reference: J. Huang, Y. Dai, L. Huang, “Adaptive Best-of-Both-Worlds Algorithm for Heavy-Tailed Multi-Armed Bandits,” ICML 2022.



清华大学 交叉信息研究院
Institute for Interdisciplinary Information Sciences, Tsinghua University

Thank you!

More info: <https://people.iis.tsinghua.edu.cn/~huang>