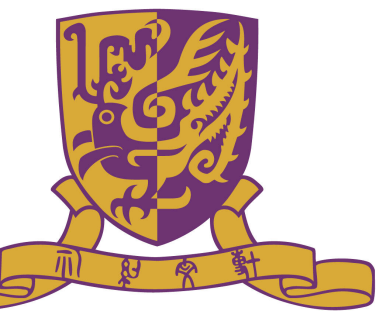


Almost Optimal Algorithms for Linear Stochastic Bandits with Heavy-Tailed Payoffs



Han Shao* Xiaotian Yu* Irwin King Michael R. Lyu

Department of Computer Science and Engineering, The Chinese University of Hong Kong, Shatin, N.T., Hong Kong
{hshao, xtyu, king, lyu}@cse.cuhk.edu.hk

Introduction

Background: Why heavy tails in bandits

- Generally, payoffs in Multi-Armed Bandits (MAB) are assumed under sub-Gaussian noises.
- In practice, many scenarios contain heavy-tailed noises, e.g., extreme payoffs in financial markets.
- Existing work has not achieved the optimal regret of linear stochastic bandits with heavy tails.

Practical motivation: An illustration

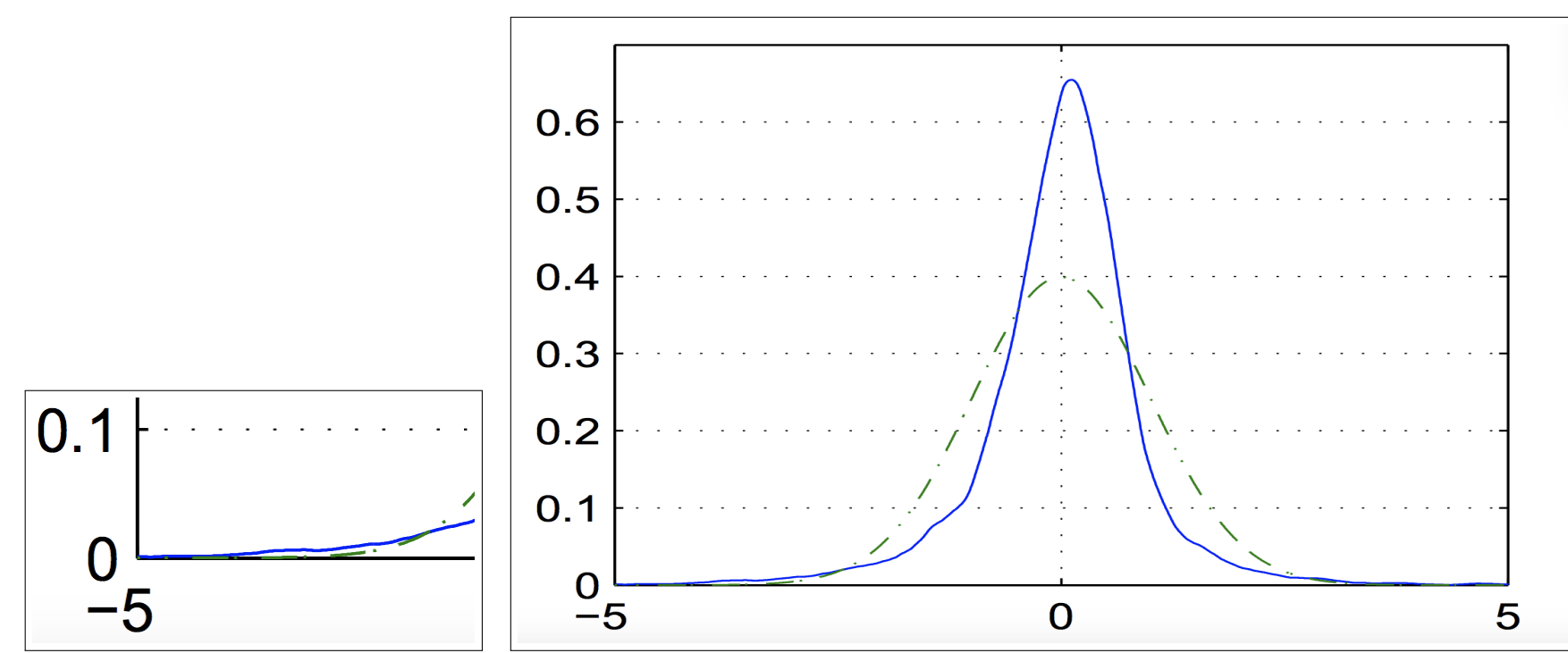


Figure 1: High-probability extreme returns in finance.

In the figure, the blue-solid line denotes the NAS-DAQ returns for the last 20 years, and the green-dashed line denotes a fitted Gaussian distribution.

The problem of LinBET

Definition 1. Given a decision set D_t for time step $t = 1, \dots, T$, an algorithm \mathcal{A} , of which the goal is to maximize cumulative payoffs over T rounds, chooses an arm $x_t \in D_t$. With \mathcal{F}_{t-1} , the observed stochastic payoff $y_t(x_t)$ is conditionally heavy-tailed, i.e., $\mathbb{E}[|y_t|^{1+\epsilon} | \mathcal{F}_{t-1}] \leq b$ or $\mathbb{E}[|y_t - \langle x_t, \theta_* \rangle|^{1+\epsilon} | \mathcal{F}_{t-1}] \leq c$, where $\epsilon \in (0, 1]$, and $b, c \in (0, +\infty)$.

Our contributions

- We provide the lower bound for the problem of LinBET as $\Omega(T^{\frac{1}{1+\epsilon}})$, where $\epsilon \in (0, 1]$.
- We develop two novel bandit algorithms, which are named as MENU and TOFU. Both algorithms achieve the regret $\tilde{O}(T^{\frac{1}{1+\epsilon}})$ with high probability.
- We conduct experiments to demonstrate the effectiveness of our proposed algorithms.

Lower Bound of LinBET

Theorem for Lower Bound. We define a set $S_d \triangleq \{(\theta_1, \dots, \theta_d) : \forall i \in [d/2], (\theta_{2i-1}, \theta_{2i}) \in \{(2\Delta, \Delta), (\Delta, 2\Delta)\}\}$ with $\Delta \in (0, 1/d]$ and $D_{(d)} \triangleq \{(x_1, \dots, x_d) \in \mathbb{R}_+^d : x_1 + x_2 = \dots = x_{d-1} + x_d = 1\}$. If θ_* is chosen uniformly at random from S_d , and the payoff for each $x \in D_{(d)}$ is in $\{0, (1/\Delta)^{\frac{1}{1+\epsilon}}\}$ with mean $\theta_*^\top x$, we have $\forall \mathcal{A}$ and $\forall T \geq (d/12)^{\frac{1}{1+\epsilon}}$,

$$\mathbb{E}[R(\mathcal{A}, T)] \geq \frac{d}{192} T^{\frac{1}{1+\epsilon}}.$$

TOFU and Upper Bound

Comparison between our TOFU and CRT

- TOFU truncates all historical payoffs in terms of $\{\mu_i\}_{i=1}^d$ for each round.
- CRT only truncates the payoff at time t and does not revisit historical information.

Theorem for TOFU. Assume that for all t and $x_t \in D_t$ with $\|x_t\|_2 \leq D$, $\|\theta_*\|_2 \leq S$, $|x_t^\top \theta_*| \leq L$ and $\mathbb{E}[|y_t|^{1+\epsilon} | \mathcal{F}_{t-1}] \leq b$. Then, with probability at least $1 - \delta$, for every $T \geq 1$, the regret of the TOFU algorithm satisfies

$$R(\text{TOFU}, T) \leq \tilde{O}\left(db^{\frac{1}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}\right).$$

Algorithm 1 TOFU

```

1: input  $d, b, \epsilon, \delta, \lambda, T, \{D_t\}_{t=1}^T$ 
2: initialization:  $V_0 = \lambda I_d, C_0 = \mathbf{B}(\mathbf{0}, S)$ 
3: for  $t = 1, 2, \dots, T$  do
4:    $b_t = \left(\frac{b}{\log(\frac{2dT}{\delta})}\right)^{\frac{1}{1+\epsilon}} t^{\frac{1-\epsilon}{2(1+\epsilon)}}$ 
5:    $(x_t, \tilde{\theta}_t) = \arg \max_{(x, \theta) \in D_t \times C_{t-1}} \langle x, \theta \rangle$   $\triangleright$  to select an arm
6:   Play  $x_t$  and observe a payoff  $y_t$ 
7:    $V_t = V_{t-1} + x_t x_t^\top$  and  $X_t^\top = [x_1, \dots, x_t]$ 
8:    $[u_1, \dots, u_d]^\top = V_t^{-1/2} X_t^\top$ 
9:   for  $i = 1, \dots, d$  do
10:     $Y_i^\dagger = (y_1 \mathbf{1}_{u_{i,1} y_1 \leq b_t}, \dots, y_t \mathbf{1}_{u_{i,t} y_t \leq b_t})$   $\triangleright$  to truncate the payoffs
11:   end for
12:    $\theta_t^\dagger = V_t^{-1/2} (u_1^\top Y_1^\dagger, \dots, u_d^\top Y_d^\dagger)$ 
13:    $\beta_t = 4\sqrt{db}^{\frac{1}{1+\epsilon}} \left(\log\left(\frac{2dT}{\delta}\right)\right)^{\frac{1}{1+\epsilon}} t^{\frac{1-\epsilon}{2(1+\epsilon)}} + \lambda^{\frac{1}{2}} S$ 
14:   Update  $C_t = \{\theta : \|\theta - \theta_t^\dagger\|_{V_t} \leq \beta_t\}$ 
15: end for
```

MENU and Upper Bound

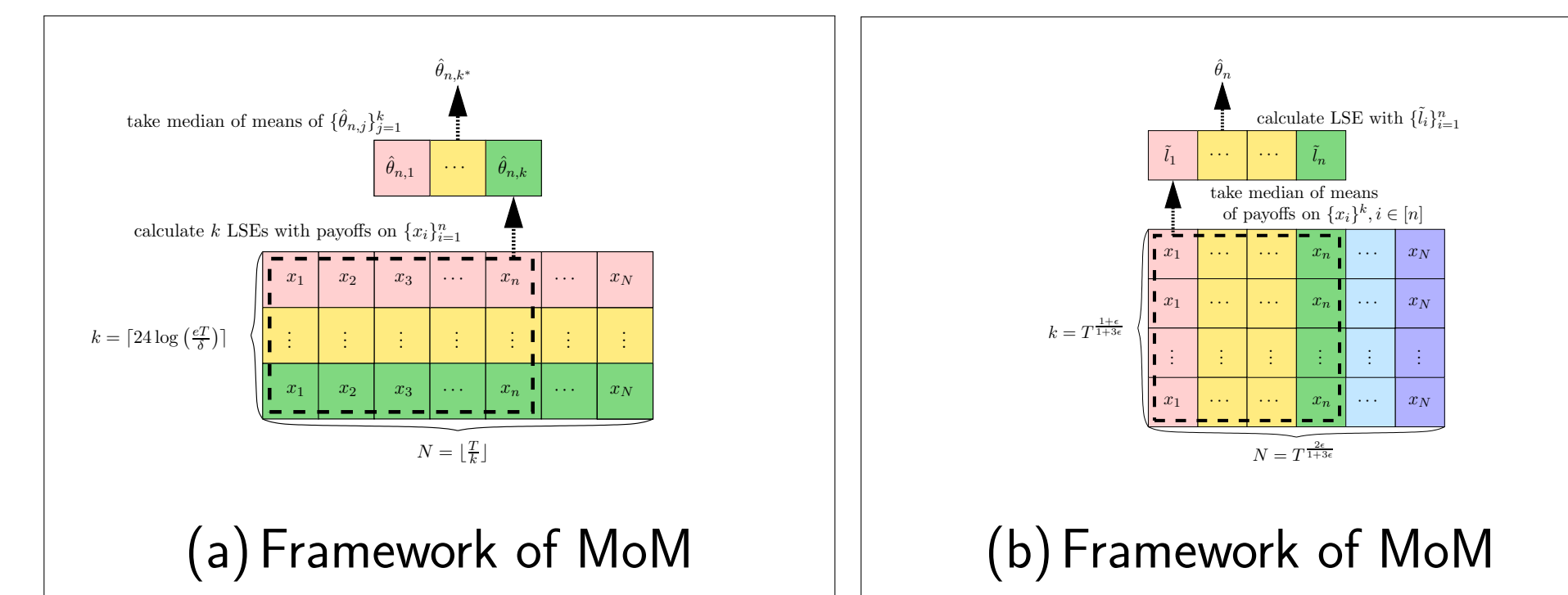


Figure 2: Comparison between our MENU and MoM.

Theorem for MENU. Assume that for all t and $x_t \in D_t$ with $\|x_t\|_2 \leq D$, $\|\theta_*\|_2 \leq S$, $|x_t^\top \theta_*| \leq L$ and $\mathbb{E}[|y_t|^{1+\epsilon} | \mathcal{F}_{t-1}] \leq c$. Then, with probability at least $1 - \delta$, for every $T \geq 256 + 24 \log(e/\delta)$, the regret of the MENU algorithm satisfies

$$R(\text{MENU}, T) \leq \tilde{O}\left(d^{\frac{1}{1+\epsilon} + \frac{1}{2}} c^{\frac{1}{1+\epsilon}} T^{\frac{1}{1+\epsilon}}\right).$$

Algorithm 2 MENU

```

1: input  $d, c, \epsilon, \delta, \lambda, S, T, \{D_n\}_{n=1}^N$ 
2: initialization:  $k = \lceil 24 \log(\frac{eT}{\delta}) \rceil, N = \lfloor \frac{T}{k} \rfloor, V_0 = \lambda I_d, C_0 = \mathbf{B}(\mathbf{0}, S)$ 
3: for  $n = 1, 2, \dots, N$  do
4:    $(x_n, \tilde{\theta}_n) = \arg \max_{(x, \theta) \in D_n \times C_{n-1}} \langle x, \theta \rangle$   $\triangleright$  to select an arm
5:   Play  $x_n$  with  $k$  times and observe payoffs  $y_{n,1}, y_{n,2}, \dots, y_{n,k}$ 
6:    $V_n = V_{n-1} + x_n x_n^\top$ 
7:   For  $j \in [k], \hat{\theta}_{n,j} = V_n^{-1} \sum_{i=1}^n y_{i,j} x_i$   $\triangleright$  to calculate LSE for the  $j$ -th group
8:   For  $j \in [k]$ , let  $r_j$  be the median of  $\{\|\hat{\theta}_{n,j} - \hat{\theta}_{n,s}\|_{V_n} : s \in [k] \setminus j\}$ 
9:    $k^* = \arg \min_{j \in [k]} r_j$   $\triangleright$  to take median of means of estimates
10:   $\beta_n = 3 \left( (9dc)^{\frac{1}{1+\epsilon}} n^{\frac{1-\epsilon}{2(1+\epsilon)}} + \lambda^{\frac{1}{2}} S \right)$ 
11:   $C_n = \{\theta : \|\theta - \hat{\theta}_{n,k^*}\|_{V_n} \leq \beta_n\}$   $\triangleright$  to update the confidence region
12: end for
```

Discussions on Theoretical Results

Future directions

- The lower bound can be further improved by the dimension d or by the moment parameter b or c .
- The regret of MENU on d can be improved.
- New algorithms achieve problem-dependent regret upper bounds with $\text{polylog}(T)$.

Experiments and Conclusion

Datasets

Table 1: Statistics of synthetic datasets in experiments. For Student's t -distribution, ν denotes the degree of freedom, l_p denotes the location, s_p denotes the scale. For Pareto distribution, α denotes the shape and s_m denotes the scale. NA denotes not available.

dataset	distribution {parameters}	$\{\epsilon, b, c\}$
S1	Student's t -distribution $\{\nu = 3, l_p = 0, s_p = 1\}$	$\{1.00, \text{NA}, 3.00\}$
S2	Student's t -distribution $\{\nu = 3, l_p = 0, s_p = 1\}$	$\{1.00, \text{NA}, 3.00\}$
S3	Pareto distribution $\{\alpha = 2, s_m = \frac{x_i^\top \theta_*}{2}\}$	$\{0.50, 7.72, \text{NA}\}$
S4	Pareto distribution $\{\alpha = 2, s_m = \frac{x_i^\top \theta_*}{2}\}$	$\{0.50, 54.37, \text{NA}\}$

Results (https://github.com/Aaronyxt/LinBET_nips2018)

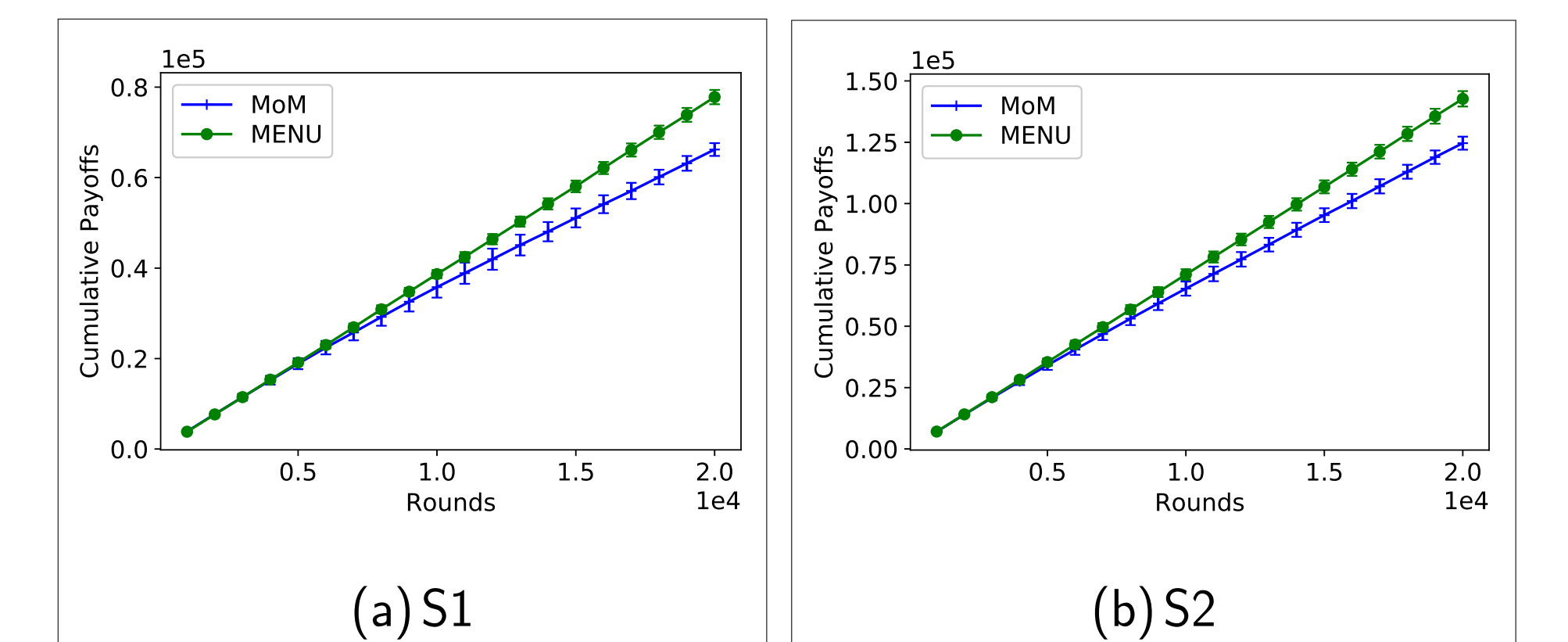


Figure 3: Comparison of cumulative payoffs for S1 and S2.

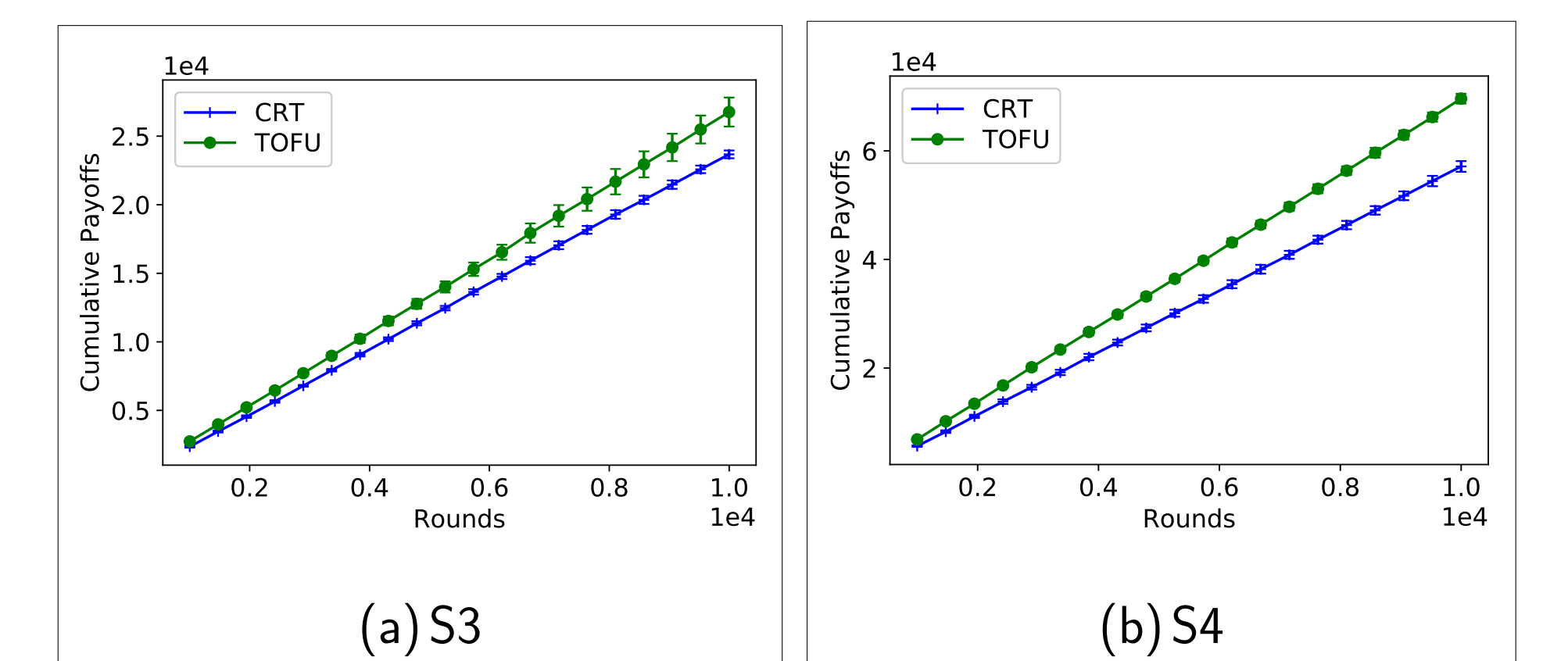


Figure 4: Comparison of cumulative payoffs for S3 and S4.

Conclusion

We broke the assumption of sub-Gaussian noises in payoffs of bandits. We rigorously analyzed the lower bound of LinBET, and developed two novel bandit algorithms with regret upper bounds matching the lower bound up to polylogarithmic factors.