

Demystifying Online Clustering of Bandits: Enhanced Exploration Under Stochastic and Smoothed Adversarial Contexts

Zhuohua Li[†], Maoli Liu[†], Xiangxiang Dai, and John C.S. Lui^{*}
The Chinese University of Hong Kong



INTRODUCTION

• We address an open problem in online clustering of bandits, an extension of contextual linear bandits that groups similar users into clusters, utilizing shared features to improve learning efficiency.

PROBLEM FORMULATION

- There are u users. Each user $i \in [u]$ is associated with an unknown preference vector $\theta_i \in \mathbb{R}^d$.
- The users are separated into m ($m \ll u$) disjoint clusters, such that:
 - Users i, j in the same cluster satisfy $\theta_i = \theta_j$.
 - Users i, j from different clusters satisfy $\|\theta_i - \theta_j\| \geq \gamma$.
- At each round $t = 1, 2, \dots, T$, the learner receives a user index $i_t \in [u]$ and a finite set of arms $\mathcal{A}_t \subset \mathcal{A} \subset \mathbb{R}^d$ where $|\mathcal{A}_t| = K$. Each arm $a \in \mathcal{A}$ is associated with a feature vector $x_a \in \mathbb{R}^d$. The learner assigns an appropriate cluster V_t for user i_t , recommends an arm $a_t \in \mathcal{A}_t$, and receives a reward $r_t = x_{a_t}^\top \theta_{i_t} + \eta_t$, where η_t is noise.
- Let $a_t^* = \arg \max_{a \in \mathcal{A}_t} x_a^\top \theta_{i_t}$ be the optimal arm at time t . The goal is to minimize the expected cumulative regret:

$$\mathbb{E}[R(T)] = \mathbb{E} \left[\sum_{t=1}^T \left(x_{a_t^*}^\top \theta_{i_t} - x_{a_t}^\top \theta_{i_t} \right) \right]$$

OPEN PROBLEM

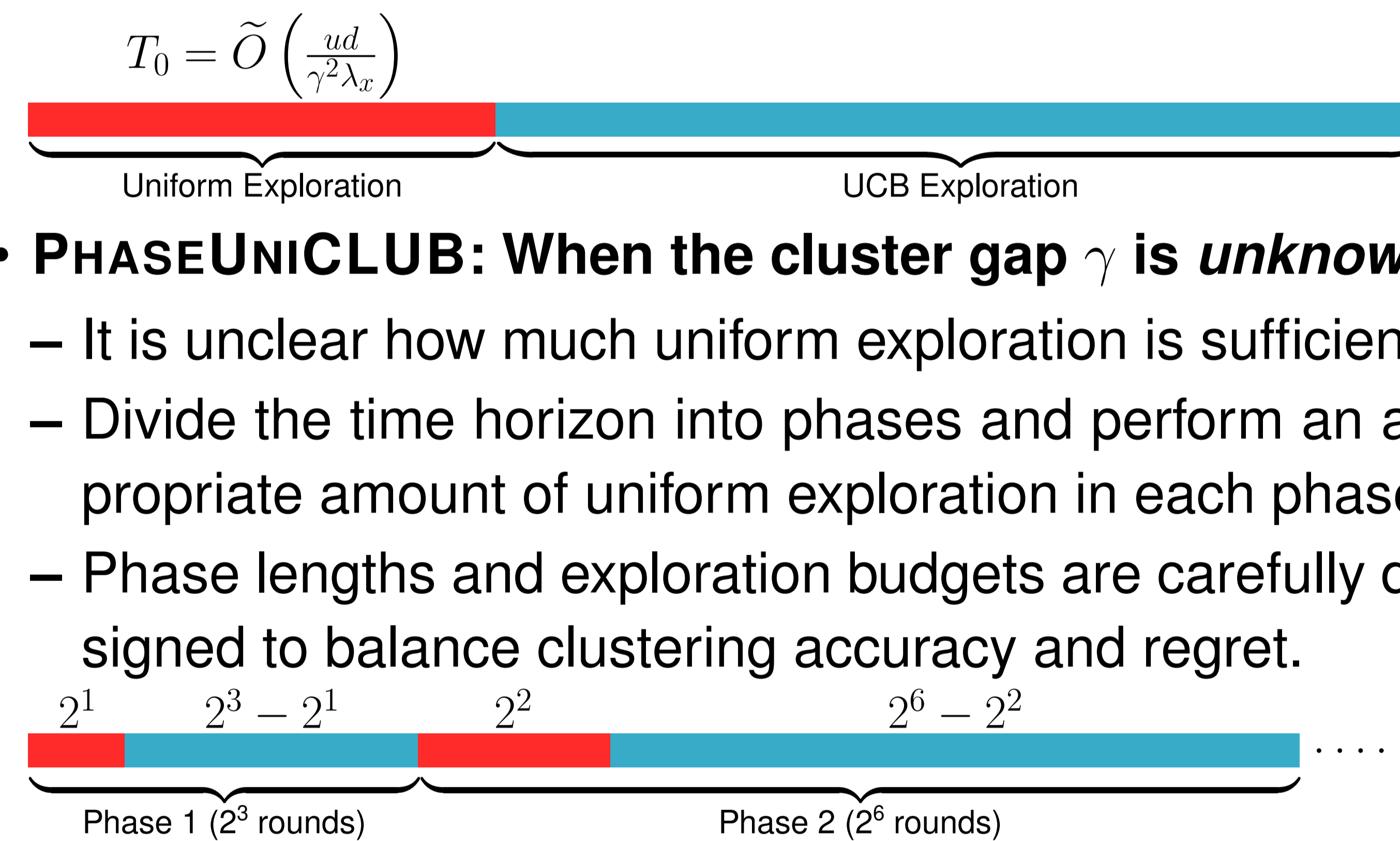
- Existing algorithms (Gentile et al., 2014) rely on strong data diversity assumptions:
 1. At each time t , vectors $\{x_a\}_{a \in \mathcal{A}_t}$ are i.i.d. sampled from a fixed distribution \mathcal{X} with $\lambda_{\min}(\mathbb{E}[\mathcal{X}\mathcal{X}^\top]) = \lambda_x$;
 2. For any unit vector $z \in \mathbb{R}^d$, $(z^\top \mathcal{X})^2$ is σ^2 -sub-Gaussian with $\sigma^2 \leq \frac{\lambda_x^2}{8 \log(4K)}$.
- Open problem posed by Gentile et al. (2014): Can we remove the i.i.d. and other statistical assumptions?
- Some follow-up work weakens these assumptions but suffers from deteriorated regret bounds.

CONTRIBUTIONS

- **Regret under Weaker Assumptions:** We solve the open problem by proposing UNICLUB and PHASEUNICLUB, which rely solely on Assumption 1 with regret $\tilde{O}(\sqrt{T})$.
- **Removal of i.i.d. Assumption:** We remove the i.i.d. assumption by introducing the *smoothed adversarial context setting* and proposing SACLUB with regret $\tilde{O}(\sqrt{T})$.

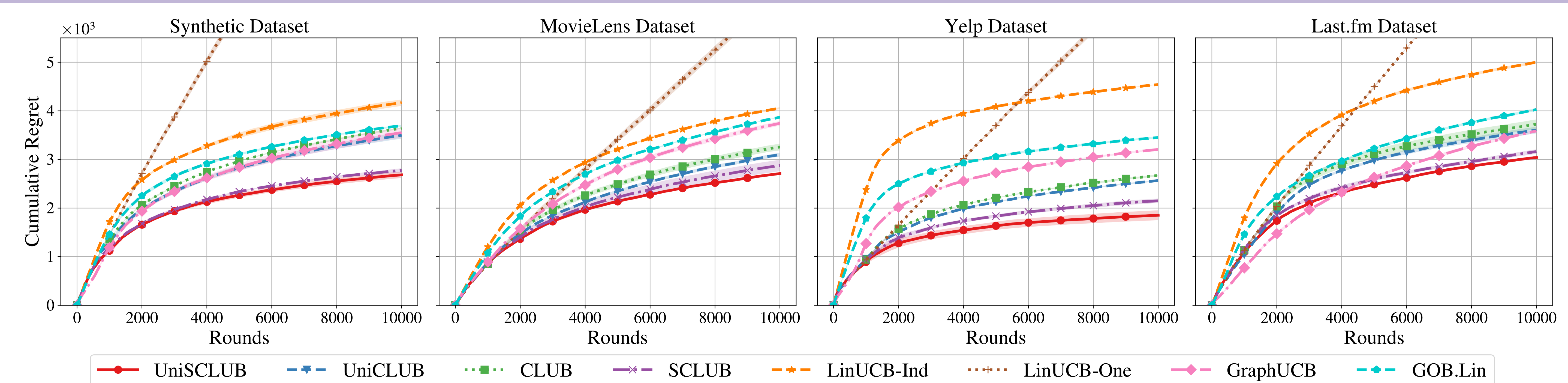
METHODOLOGIES

- **Key Idea: Add additional uniform exploration to ensure accurate clustering while maintaining low regret.**
- **UNICLUB: When the cluster gap γ is known**
 - Use γ to determine when clustering is reliable.
 - Uniform exploration until all θ_i are estimated accurately.



- **PHASEUNICLUB: When the cluster gap γ is unknown**
 - It is unclear how much uniform exploration is sufficient.
 - Divide the time horizon into phases and perform an appropriate amount of uniform exploration in each phase.
 - Phase lengths and exploration budgets are carefully designed to balance clustering accuracy and regret.

EVALUATIONS



SMOOTHED ADVERSARIAL CONTEXT SETTING

- **Key Idea: The intrinsic diversity of contexts makes explicit exploration unnecessary.**
- Each feature vector is first arbitrarily chosen by an adversary and then perturbed by a noise vector sampled from a truncated multivariate Gaussian distribution.
- This setup is more practical and aligns more closely with the original setting of contextual linear bandits.

THEORETICAL RESULTS

Theorem 1. With Assumption 1 and assuming the cluster gap γ is known, algorithm UNICLUB satisfies:

$$\mathbb{E}[R(T)] = \tilde{O}\left(\frac{ud}{\gamma^2 \lambda_x} + d\sqrt{mT}\right).$$

Theorem 2. With Assumption 1, PHASEUNICLUB satisfies:

$$\mathbb{E}[R(T)] = \tilde{O}\left(\frac{ud}{\gamma^5 \lambda_x^2} + \left(\frac{ud}{\lambda_x}\right)^{\frac{2}{3}} T^{\frac{1}{3}} + d\sqrt{mT}\right).$$

Theorem 3. Under the smoothed adversarial context setting, algorithm SACLUB satisfies:

$$\mathbb{E}[R(T)] = \tilde{O}\left(\frac{ud}{\gamma^2 \tilde{\lambda}_x} + d\sqrt{mT}\right),$$

where $\tilde{\lambda}_x = \frac{C}{\log K}$ for some constant C .