

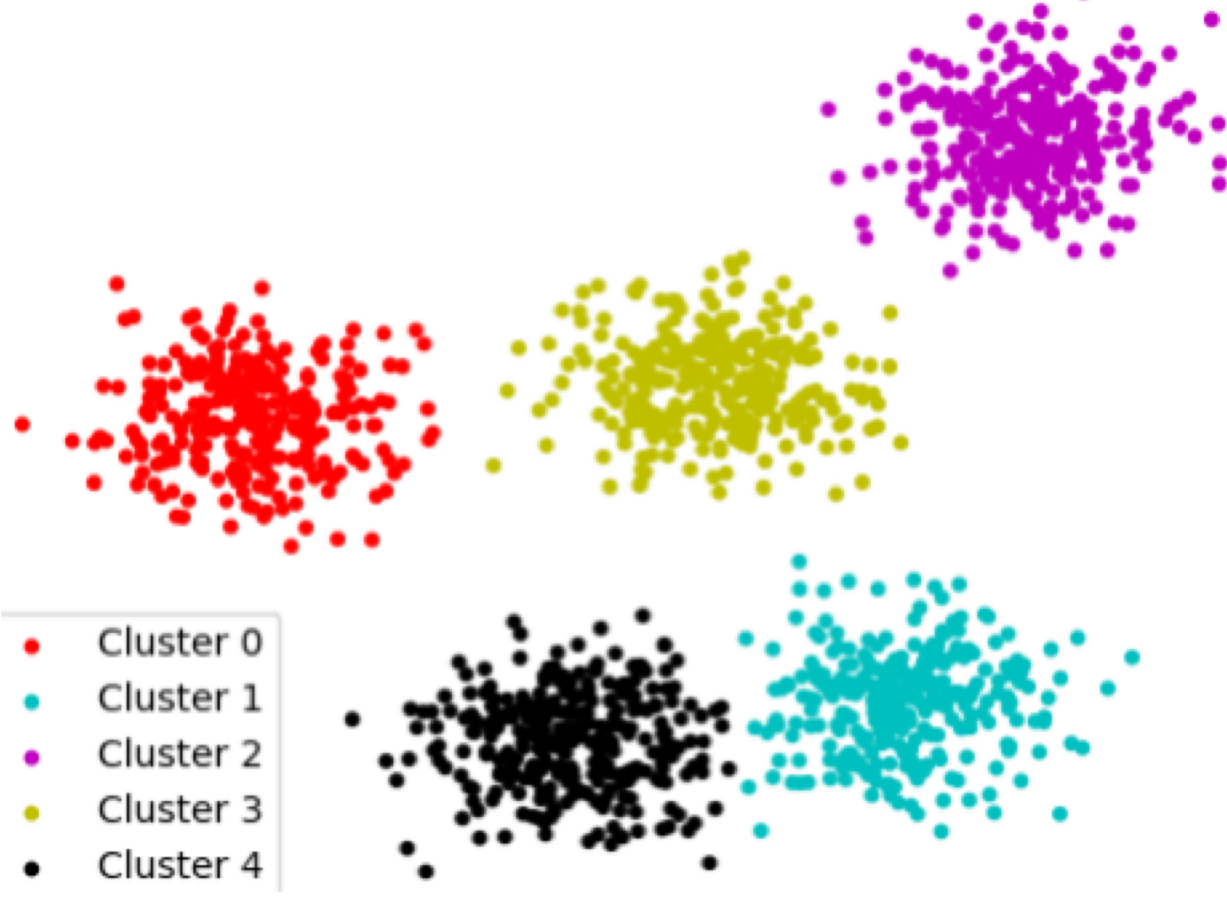
# Online Clustering of Contextual Cascading Bandits

Shuai Li<sup>\*1</sup> Shengyu Zhang<sup>\*1,2</sup>

1. The Chinese University of Hong Kong. 2. Tencent

## Motivation

- Cascading feedback
  - Scenarios: web search results, online recommendation systems, ...
  - Model: On an ordered list
    - A user goes through the list from top down,
    - stops at the first satisfactory item, and
    - clicks it.
  - Task: Use this online feedback to help improve future list recommendation
- Contexts: Features for an item or and item-user pair
  - Important for recommendations
- Combinatorial:
  - An action is an ordered sequence of items
- Clustering
  - Users have a clustering structure
  - We only see the user indices
  - Need to learn the user similarities online



## Algorithm: CLUB-cascade

- Parameters:  $\lambda, \alpha, \beta > 0$
  - Initialization:
    - $G$  is a complete graph over users;
    - $S_i = 0_{d \times d}, b_i = 0_{d \times 1}, T_i = 0$  for any user  $i$ .
  - For all  $t = 1, 2, \dots, n$  do
    - Obtain user index  $I_t$  and item set  $D_t \subset \mathbb{R}^{d \times 1}$
    - Find the connected component  $V_t$  of user  $I_t$  and compute
 
$$M = \lambda I + \sum_{i \in V_t} S_i, b = \sum_{i \in V_t} b_i, T = \sum_{i \in V_t} T_i$$
    - Compute
 
$$\hat{\theta} = M^{-1}b$$

Linear regression
    - For any  $x \in D_t$ , compute
 
$$U_t(x) = \min\{\hat{\theta}^\top x + \beta \|x\|_{M^{-1}}, 1\}$$

exploitation

exploration
    - Recommend the  $K$  items with largest  $U_t$  values and observe  $K_t$ ;  $w_t(x_k^t), k \leq K_t$ .
    - Update statistics
 
$$S_i = S_i + \sum_{k=1}^{K_t} x_{t,k} x_{t,k}^\top,$$

$$b_i = b_i + x_{t,K_t} w_t(x_{t,K_t}),$$

$$T_i = T_i + K_t$$

$$\hat{\theta}_i = (\lambda I + S_i)^{-1} b_i$$
    - Delete edge  $(i, l)$  if
 
$$\|\hat{\theta}_i - \hat{\theta}_l\| \geq \alpha \left( \frac{\sqrt{\ln T_i}}{T_i} + \frac{\sqrt{\ln T_l}}{T_l} \right)$$
- End for  $t$

## Setting

- $n_u$  of users.
- Each action is an ordered list of  $K$  items.
- At time step  $t$ ,
  - User  $I_t$  comes to be served with items  $D_t \subset \mathbb{R}^d$ .
  - Let  $\mathcal{H}_t$  denote the history so far.
  - The learning agent recommends  $A_t = (x_1^t, \dots, x_K^t)$  to the user.
  - The user checks from the first item of  $A_t$  and stops at  $K_t$ -th item.
  - The learning agent observes the weights of first  $K_t$  base arms in  $A_t$ ,  $w_t(a_k^t), k \leq K_t$ .
- Assume that given  $\mathcal{H}_t, w_t(a)$ 's are mutually independent Bernoulli random variables with
 
$$\mathbb{E}[w_t(a) | \mathcal{H}_t] = \theta_{I_t}^\top x_{t,a}$$
 for some unknown  $\theta_{I_t} \in \mathbb{R}^{d \times 1}$  with  $\|\theta_{I_t}\|_2 \leq 1, 0 \leq \theta_{I_t}^\top x_{t,a} \leq 1$ .
- Cluster regularity: All users in the same cluster have the same  $\theta$ . Users in different clusters have noticeably different  $\theta$ 's:
 
$$\|\theta - \theta'\| \geq \gamma > 0.$$
- User uniformness: At each time, the user is drawn uniformly from the set of all users, independently over the past.
- Item regularity: At each time step, the items are drawn independently from a fixed distribution where  $\mathbb{E}[xx^\top]$  has minimal eigenvalue  $\lambda_x > 0$ .
- The regret of action  $A$  on time  $t$  is
 
$$R(t, A) = f_t^* - f(A, w_t)$$
 where  $f_t^* = \max_{A^*} f(A^*, w_t)$
- Task: Minimize the cumulative regret of  $n$  rounds
 
$$R(n) = \mathbb{E} \left[ \sum_{t=1}^n R(t, A_t) \right].$$

## Theoretical analysis

**Theorem 1.** Let  $\beta = \sqrt{d \ln(1 + n/\lambda d) + 2 \ln 4mn} + \sqrt{\lambda}$  and  $\alpha = 4\sqrt{d}/\lambda_x$ . Then the regret of our algorithm, CLUB-cascade, satisfies

$$R(n) = O(d\sqrt{mnK} \ln n).$$

**Corollary 2.** When the number of clusters  $m = 1$ , the regret satisfies

$$R(n) = O(d\sqrt{nK} \ln n)$$

which is better than the results in [2][3].

**Theorem 3.** Consider a generalized linear reward function

$$\mu(\theta_{I_t}^\top x_{t,a}),$$

where  $\mu$  is strictly increasing, continuously differentiable, and Lipschitz with constant  $\kappa$ . Let  $c = \inf_{a \in [-2, 2]} \mu'(a)$ . Then the regret satisfies

$$R(n) = O\left(\frac{\kappa d}{c} \sqrt{mnK} \ln n\right)$$

## Experiments

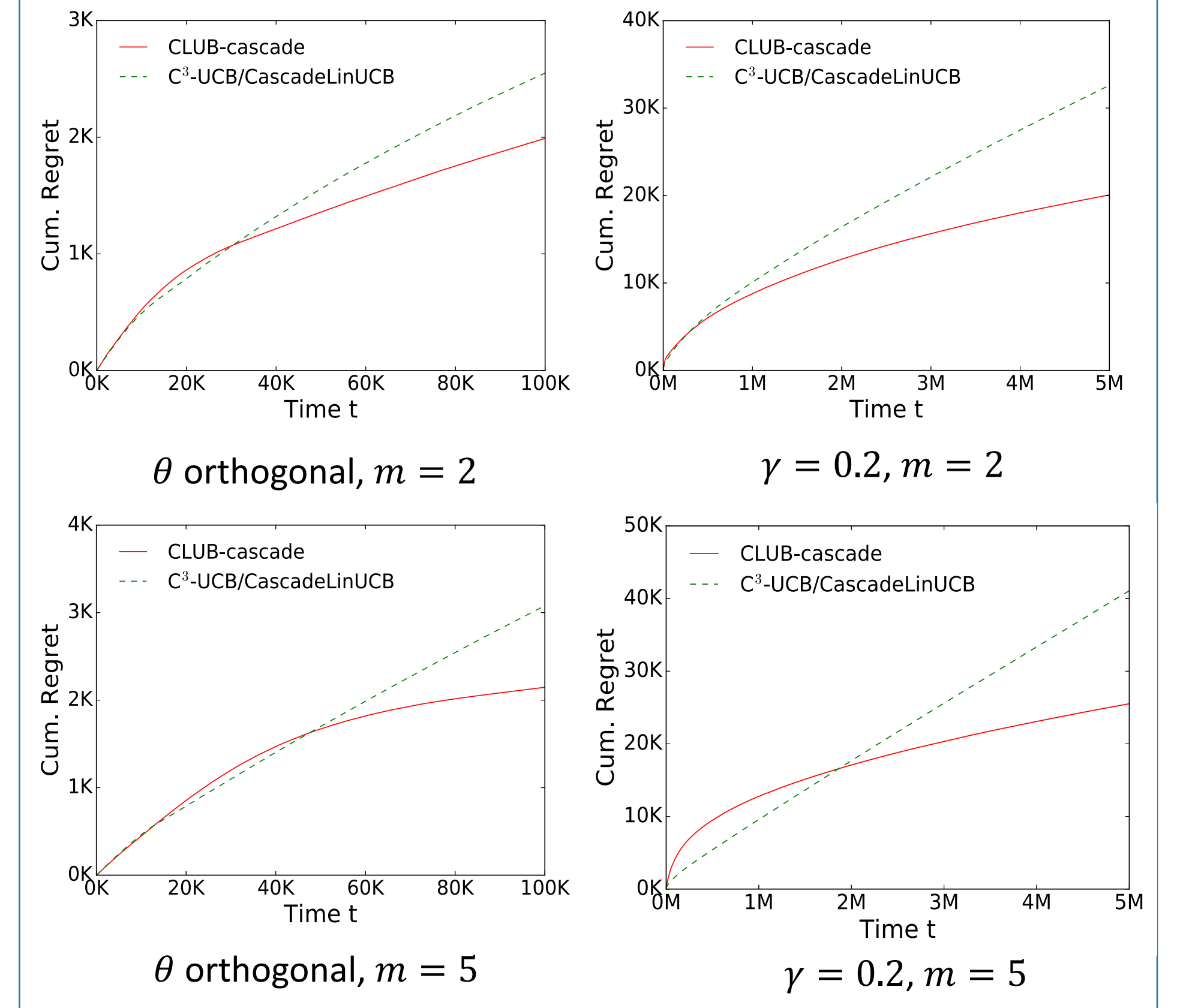


Figure 1. Experimental results for synthetic data. 40 users, 200 items,  $K = 4, d = 20$ .

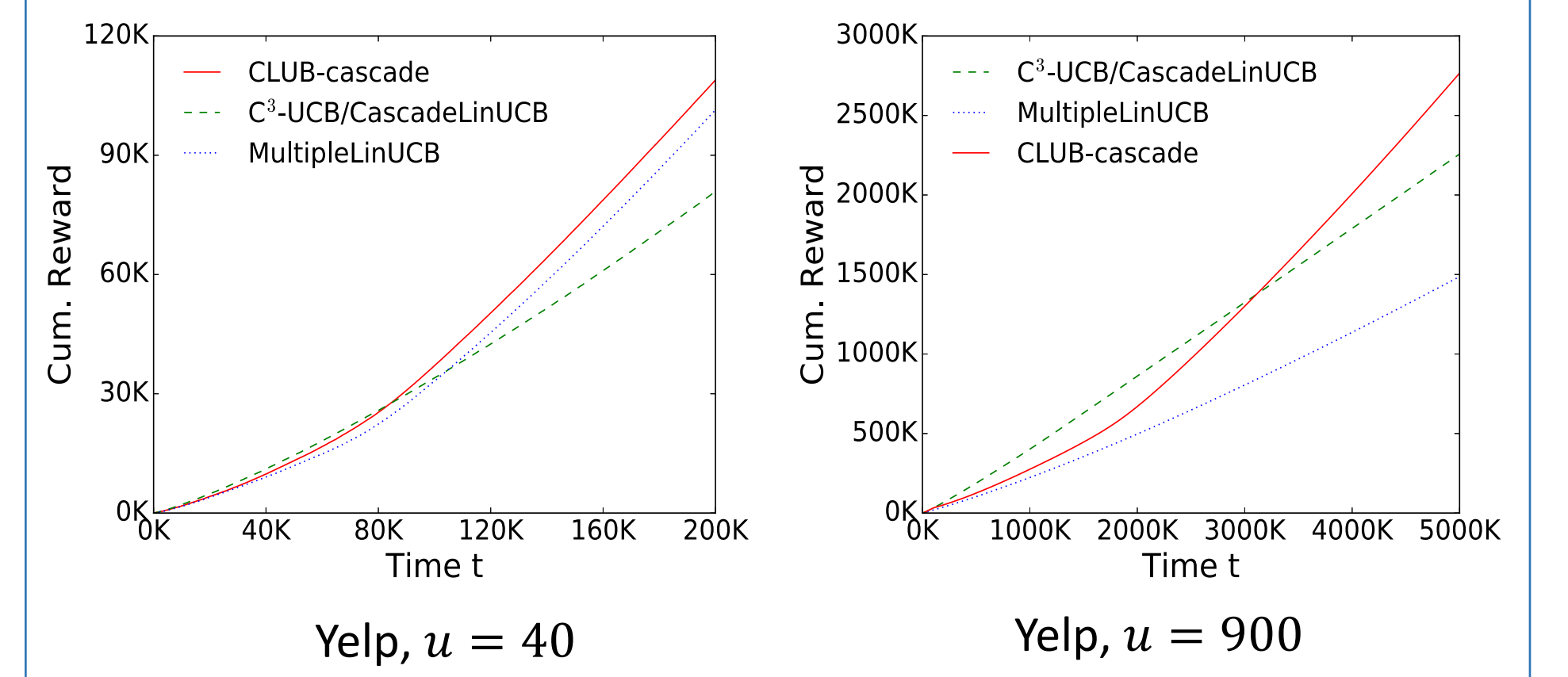


Figure 2. Cumulative rewards on Yelp dataset.  $K = 4, d = 20, L = 1k$ .

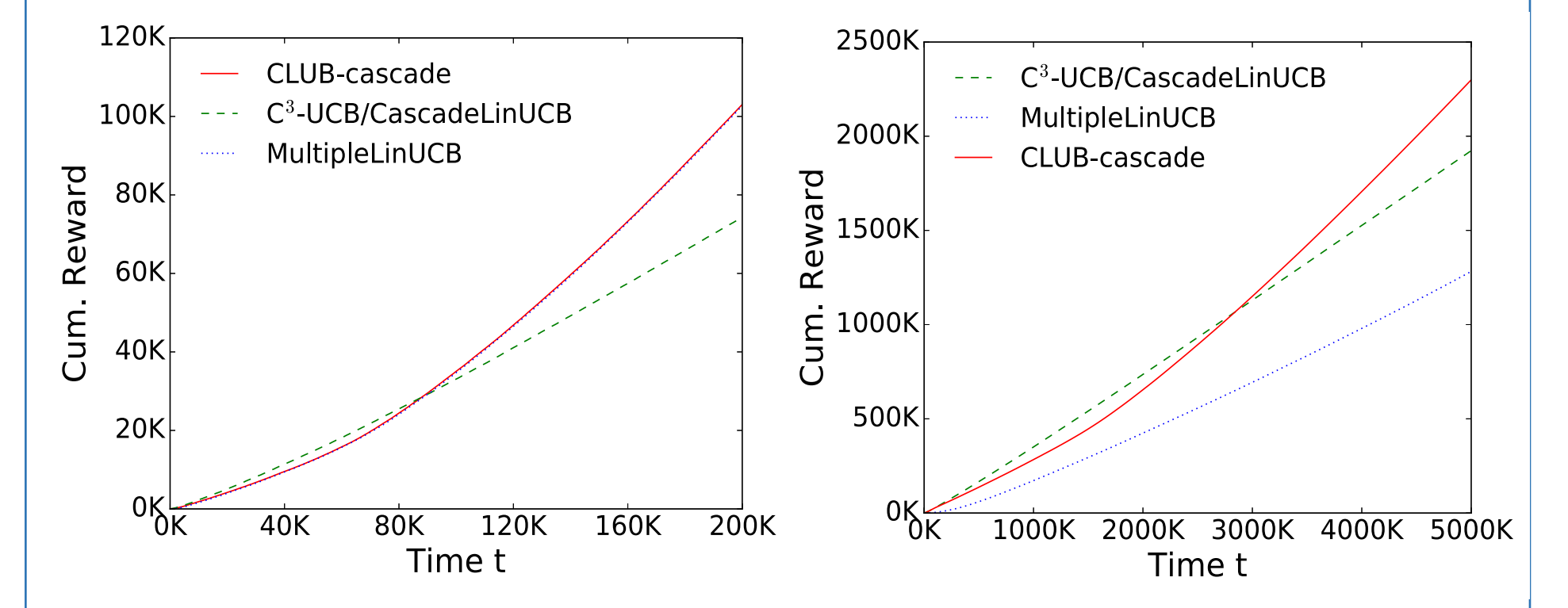


Figure 3. Cumulative rewards on MovieLens dataset.  $K = 4, d = 20, L = 1k$ .

## Conclusions

- Formulate Online Clustering of Contextual Cascading Bandits problem.
- Propose a CLUB-cascade algorithm that can learn clustering over users and, at the same time, effectively handle
  - contextual information
  - cascading feedback
- Theoretical analysis
- Empirical evaluation

## Contact

Shuai Li  
Email: shuaili@cse.cuhk.edu.hk

Shengyu Zhang  
Email: syzhang@cse.cuhk.edu.hk

## References

- Gentile, Li, and Zappella. "Online clustering of bandits." Proceedings of the 31st International Conference on Machine Learning. 2014.
- Li, Wang, Zhang, Chen. "Contextual combinatorial cascading bandits." International Conference on Machine Learning. 2016.
- Zong, Ni, Sung, Ke, Wen, Kveton. "Cascading bandits for large-scale recommendation problems." Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence. AUAI Press, 2016.