



INTRODUCTION

- Recommender systems need to learn and adapt to users' preferences by collecting their feedback, such as click rates.
- Recommender systems can also proactively query users with some questions to obtain additional feedback.

MOTIVATING EXAMPLE AND MAIN IDEAS

Which response do you prefer?
Your choice will help make ChatGPT better.

Response 1

LinUCB is implemented as a class:

```
class LinUCB:
    def __init__(self, d, alpha):
        self.alpha = alpha
        self.A = np.eye(d)
        self.b = np.zeros(d)
    def select_action(self, context):
        ...
    def update(self, action, reward):
        ...
```

Response 2

LinUCB is implemented as a function:

```
def linucb(alpha, contexts, rewards,
           actions, n_rounds):
    d = contexts.shape[1]
    A = np.eye(d)
    b = np.zeros(d)
    for t in range(n_rounds):
        ...
    return theta
```

Fig. 1: An example of conversational recommendation.

Main Idea: Query users with questions (i.e., key terms) that are closest to the direction with the maximum uncertainty.

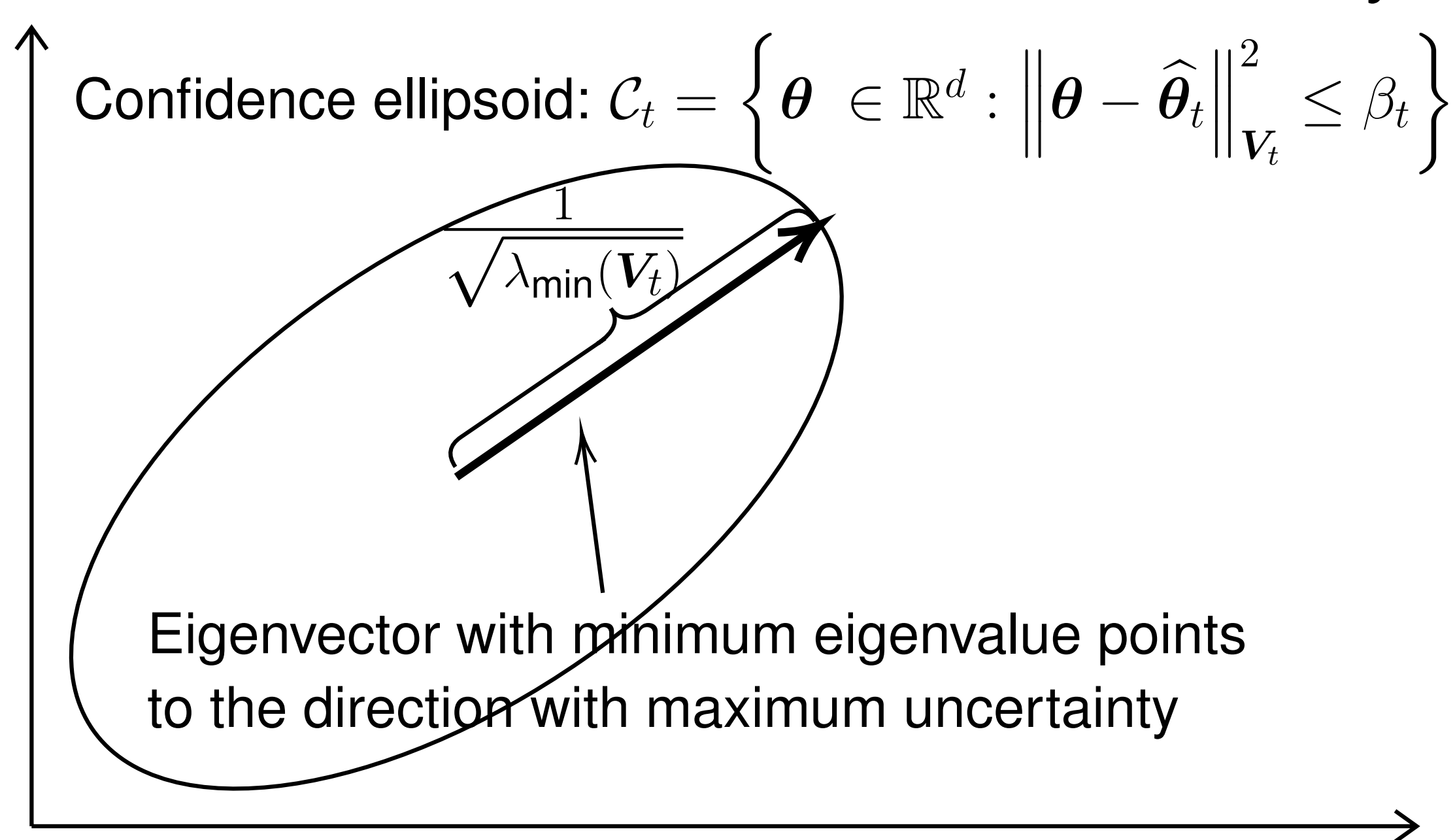


Fig. 2: Illustration of the key term selection strategy.

CONTRIBUTIONS

- Propose FedConPE, an algorithm for federated conversational bandits with heterogeneous arm sets.
- Exploit the properties of the conversational setting to enhance efficiency in both computation and communication.
- Theoretical results show that FedConPE is near-optimal.
- Experiments demonstrate that FedConPE achieves lower cumulative regret and uses fewer conversations.

PROBLEM FORMULATION

- M clients and one central server.
- Each client i has a finite set of recommended items (i.e., arms) $\mathcal{A}_i \subset \mathbb{R}^d$ and $|\mathcal{A}_i| = K$.
- The server has a finite set of questions (i.e., key terms) $\mathcal{K} \subset \mathbb{R}^d$, and we assume that \mathcal{K} is sufficiently rich.

Interaction Protocol

- For each time step $t = 1, 2, \dots, T$:
 - Each client i chooses an arm $\mathbf{a}_{i,t} \in \mathcal{A}_i$, and observes an arm-level reward $x_{i,t} = \mathbf{a}_{i,t}^\top \boldsymbol{\theta}^* + \eta_{i,t}$.
 - The server decides whether to initiate queries $\mathbf{k}_{i,t} \in \mathcal{K}$ to each client i . If so, the client observes a key term-level reward $\tilde{x}_{i,t} = \mathbf{k}_{i,t}^\top \tilde{\boldsymbol{\theta}}^* + \tilde{\eta}_{i,t}$.
- The objective is to minimize the cumulative regret:

$$R_M(T) = \sum_{i=1}^M \sum_{t=1}^T \left(\max_{\mathbf{a} \in \mathcal{A}_i} \mathbf{a}^\top \boldsymbol{\theta}^* - \mathbf{a}_{i,t}^\top \boldsymbol{\theta}^* \right),$$

while using as few queries as possible.

THEORETICAL RESULTS

Theorem 1 (Regret upper bound). *With probability at least $1 - \delta$, the cumulative regret scales as $\mathcal{O}\left(\sqrt{dMT \log \frac{KM \log T}{\delta}}\right)$.*

Theorem 2 (Regret lower bound). *For any policy that chooses at most one key term at each time step, there exists an instance of the federated conversational bandits such that the policy suffers an expected regret of at least $\Omega(\sqrt{dMT})$.*

Theorem 3 (Communication cost). *The number of scalars transmitted between the server and clients scales as $\mathcal{O}(d^2 M \log T)$.*

Theorem 4 (Conversation frequency upper bound). *The fraction of conversations relative to the total time horizon is at most $\frac{1}{NC^2}$.*

EVALUATIONS

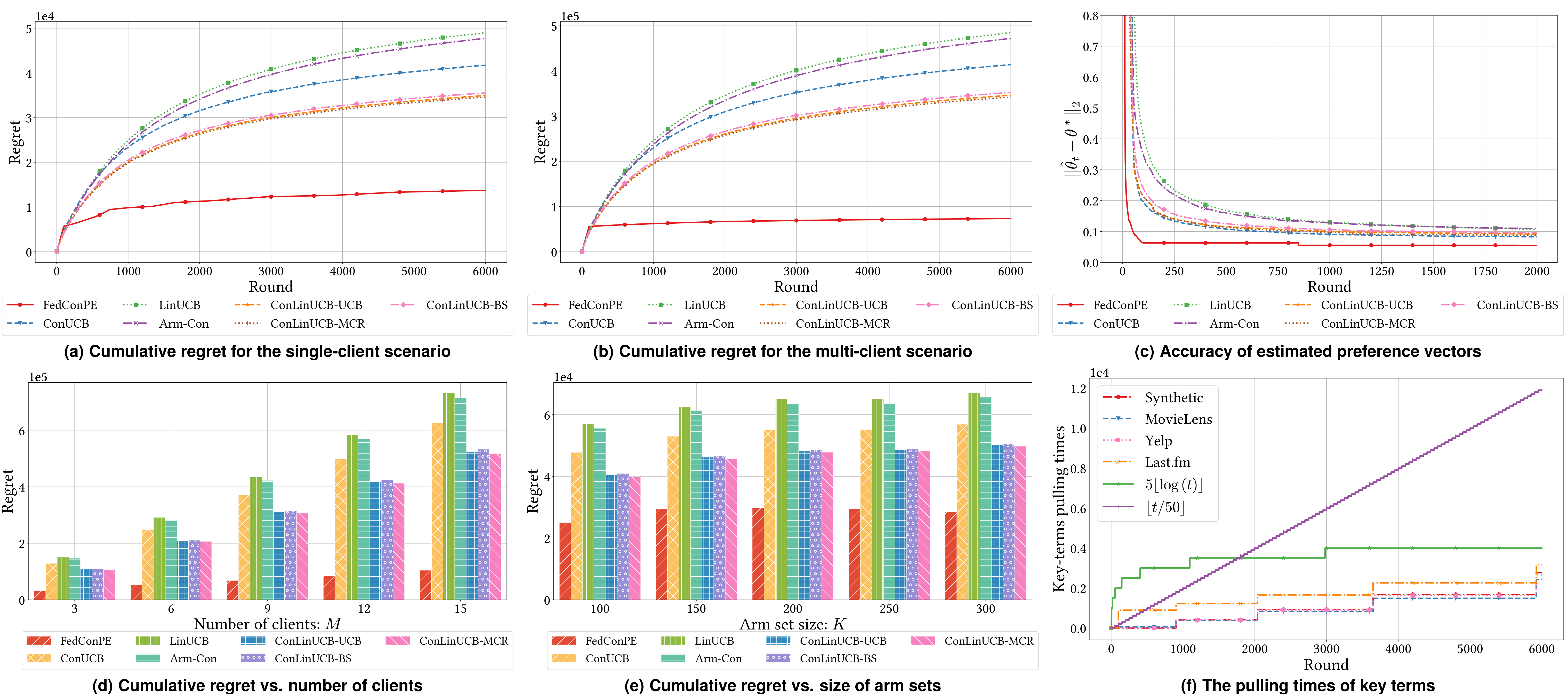


Fig. 3: Performance evaluation of FedConPE and comparison with baseline algorithms.