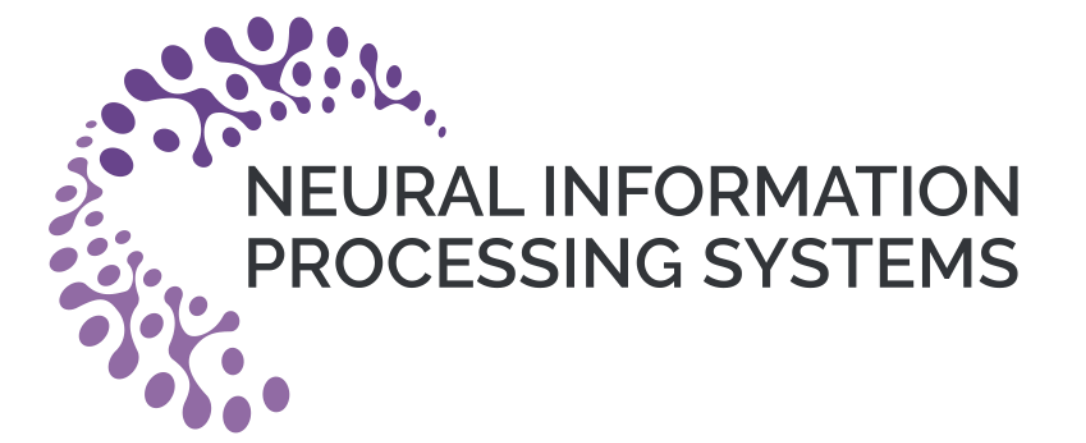


Communication Efficient Distributed Learning for Kernelized Contextual Bandits

Chuanhao Li¹, Huazheng Wang², Mengdi Wang³, and Hongning Wang¹

1. University of Virginia, 2. Oregon State University, 3. Princeton University



Abstract

We tackle the communication efficiency challenge of learning kernelized contextual bandits in a distributed setting. Despite the recent advances in communication-efficient distributed bandit learning, existing solutions are restricted to simple models like multi-armed bandits and linear bandits, which hamper their practical utility. In this paper, instead of assuming the existence of a linear reward mapping from the features to the expected rewards, we consider non-linear reward mappings, by letting agents collaboratively search in a reproducing kernel Hilbert space (RKHS). This introduces significant challenges in communication efficiency as distributed kernel learning requires the transfer of raw data, leading to a communication cost that grows linearly w.r.t. time horizon T . We address this issue by equipping all agents to communicate via a common Nystrom embedding that gets updated adaptively as more data points are collected. We rigorously proved that our algorithm can attain sub-linear rate in both regret and communication cost.

Distributed Bandit Learning

For each round $l = 1, 2, \dots, T$

For client $i = 1, 2, \dots, N$

index $t := N(l-1) + i$

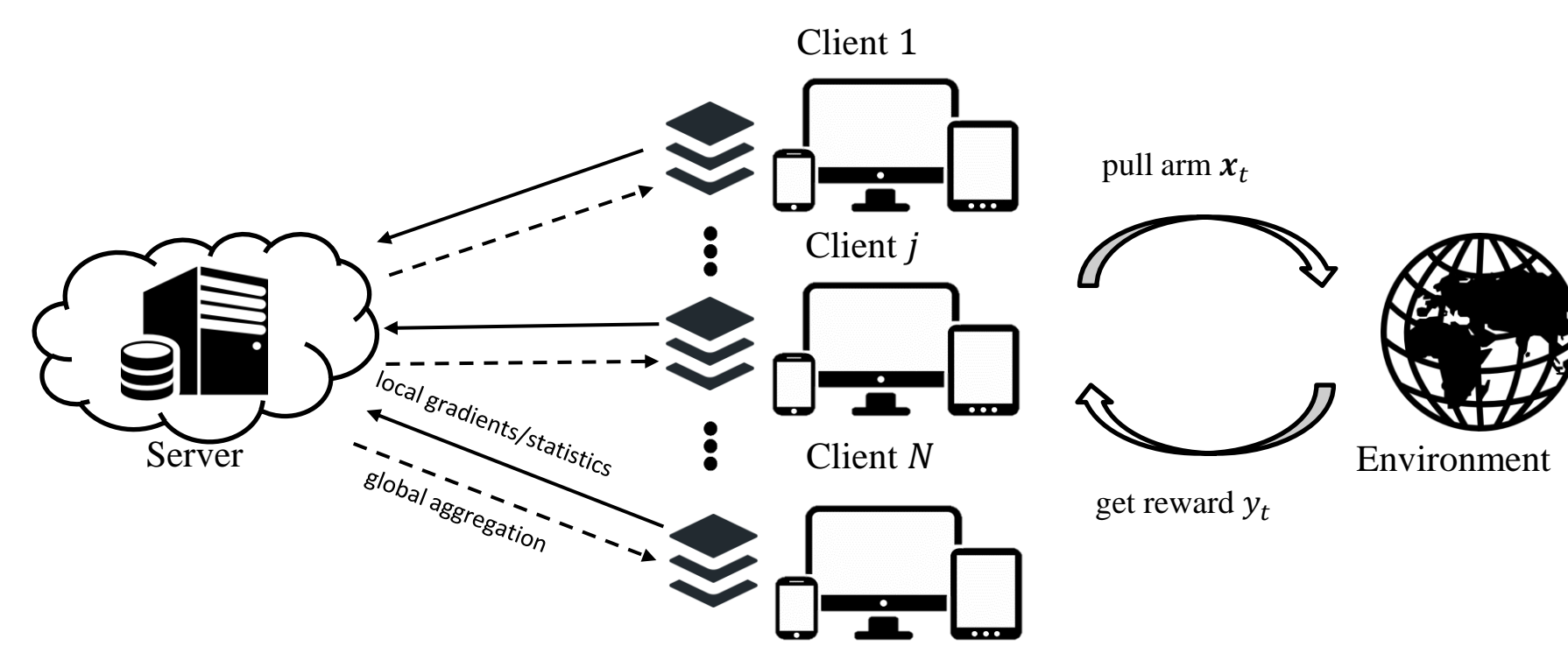
- Client i_t picks arm x_t from set \mathcal{A}_t and observes reward $y_t = f(x_t) + \eta_t$
- Communication between the server and clients

Regret & Communication

- $R_T = \sum_{t=1}^{NT} r_t$, where $r_t = \max_{x \in \mathcal{A}_t} f(x) - f(x_t)$
- C_T : total number of real numbers transferred in the system

Goal

- Incur sub-linear C_T w.r.t. T , while having near-optimal $R_T = \tilde{O}(\sqrt{NT})$



A network with N clients sequentially taking actions and receiving feedback from the environment, and a server that coordinates the communication among the clients.

Extension to Kernelized Contextual Bandits

Prior works in linear bandits & challenges in extension to kernelized bandits

Distributed linear bandits

- Joint model estimation $\hat{\theta} = \mathbf{A}^{-1}\mathbf{b}$
- Communicate local updates of $\mathbf{A} = \lambda\mathbf{I} + \mathbf{X}^T\mathbf{X} \in \mathbb{R}^{d \times d}$, $\mathbf{b} = \mathbf{X}^T\mathbf{y} \in \mathbb{R}^d$

	Regret R_T	Communication C_T
[Wang et al., ICLR' 20]	$O(d\sqrt{NT} \log NT)$	$\tilde{O}(d^3 N^{1.5})$
[Li and Wang, AISTATS' 22, He et al., NeurIPS' 22]	$O(d\sqrt{NT} \log NT)$	$\tilde{O}(d^3 N^2)$

Distributed kernelized contextual bandits:

- Assume f in RKHS: $f(x) = \phi(x)^T \theta_*$
 $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^p$ is a known feature map p is possibly infinite
 $\theta_* \in \mathbb{R}^p$ is the unknown parameter
- Search for unknown reward function f in RKHS, which is a powerful tool for optimizing black box functions

Challenge: joint kernel estimation is communication expensive

- Empirical mean and variance

$$\hat{\mu}_{t,i}(\mathbf{x}) = \mathbf{K}_{\mathcal{D}_t(i)}(\mathbf{x})^T (\mathbf{K}_{\mathcal{D}_t(i), \mathcal{D}_t(i)} + \lambda\mathbf{I})^{-1} \mathbf{y}_{\mathcal{D}_t(i)}$$

$$\hat{\sigma}_{t,i}(\mathbf{x}) = \lambda^{-1/2} \sqrt{k(\mathbf{x}, \mathbf{x}) - \mathbf{K}_{\mathcal{D}_t(i)}(\mathbf{x})^T (\mathbf{K}_{\mathcal{D}_t(i), \mathcal{D}_t(i)} + \lambda\mathbf{I})^{-1} \mathbf{K}_{\mathcal{D}_t(i)}(\mathbf{x})}$$

where

$$\mathbf{K}_{\mathcal{D}_t(i)}(\mathbf{x}) = \Phi_{\mathcal{D}_t(i)} \phi(\mathbf{x}) = [k(\mathbf{x}_s, \mathbf{x})]_{s \in \mathcal{D}_t(i)}^T \in \mathbb{R}^{|\mathcal{D}_t(i)|}$$

$$\mathbf{K}_{\mathcal{D}_t(i), \mathcal{D}_t(i)} = \Phi_{\mathcal{D}_t(i)}^T \Phi_{\mathcal{D}_t(i)} = [k(\mathbf{x}_s, \mathbf{x}_{s'})]_{s, s' \in \mathcal{D}_t(i)} \in \mathbb{R}^{|\mathcal{D}_t(i)| \times |\mathcal{D}_t(i)|}$$

Proposed Solution

Nystrom Approximation

- Approximated mean and variance

$$\tilde{\mu}_{t,i}(\mathbf{x}) = z(\mathbf{x}; \mathcal{S})^T (\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}}^T \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}} + \lambda\mathbf{I})^{-1} \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}}^T \mathbf{y}_{\mathcal{D}_t(i)})$$

$$\tilde{\sigma}_{t,i}(\mathbf{x}) = \lambda^{-1/2} \sqrt{k(\mathbf{x}, \mathbf{x}) - z(\mathbf{x}; \mathcal{S})^T \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}}^T \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}} \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}}^T \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}} + \lambda\mathbf{I}]^{-1} z(\mathbf{x}; \mathcal{S})}$$

where

$\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}} \in \mathbb{R}^{|\mathcal{D}_t(i)| \times |\mathcal{S}|}$ is obtained by applying embedding function $z(\cdot)$ to $\Phi_{\mathcal{D}_t(i)}$ embedding function $z(\cdot)$ is shared by all N clients

Event-triggered communication

- if $\sum_{s \in \mathcal{D}_t(i) \setminus \mathcal{D}_{\text{last}}(i)} \hat{\sigma}_{\text{last}, i}^2(\mathbf{x}_s) > D$ for any client i

update the shared embedding function $z(\cdot)$
synchronize embedded statistics of all clients

Algorithm 2 Approximated Distributed Kernel UCB (Approx-DisKernelUCB)

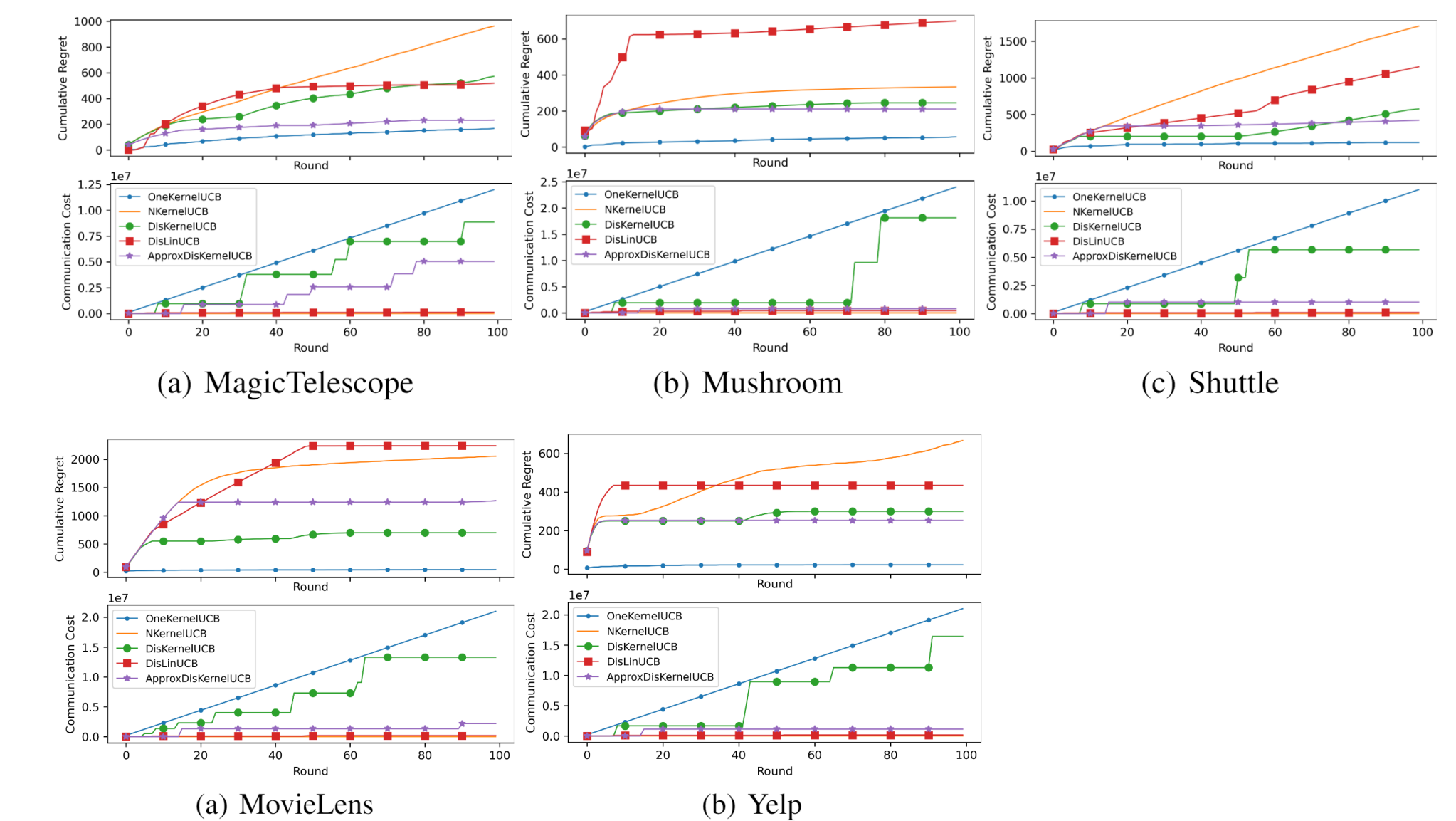
- Input:** threshold $D > 0$, regularization parameter $\lambda > 0$, $\delta \in (0, 1)$ and kernel function $k(\cdot, \cdot)$.
- Initialize** $\tilde{\mu}_{0,i}(\mathbf{x}) = 0, \tilde{\sigma}_{0,i}(\mathbf{x}) = \sqrt{k(\mathbf{x}, \mathbf{x})}, \mathcal{N}_0(i) = \mathcal{D}_0(i) = \emptyset, \forall i \in [N]; \mathcal{S}_0 = \emptyset, t_{\text{last}} = 0$
- for** round $l = 1, 2, \dots, T$ **do**
- for** client $i = 1, 2, \dots, N$ **do**
- [Client i] selects arm $\mathbf{x}_t \in \mathcal{A}_t$ according to Eq (3) and observes reward y_t , where $t := N(l-1) + i$
- [Client i] updates $\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}_t}$, $\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}_{t_{\text{last}}}}$ and $\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}_{t_{\text{last}}}}^T \mathbf{y}_{\mathcal{D}_t(i)}$ using $(z(\mathbf{x}_t; \mathcal{S}_{t_{\text{last}}}), y_t)$; sets $\mathcal{N}_t(i) = \mathcal{N}_{t-1}(i) \cup \{t\}$, and $\mathcal{D}_t(i) = \mathcal{D}_{t-1}(i) \cup \{t\}$
// Global Synchronization
- if** the event $\mathcal{U}_t(D)$ defined in Eq (4) is true **then**
- [Clients $\forall i$] sample $\mathcal{S}_t = \text{RLS}(\mathcal{N}_t(i), \bar{q}_t, \hat{\sigma}_{t_{\text{last}, i}}^2)$, and send $\{(\mathbf{x}_s, y_s)\}_{s \in \mathcal{S}_t}$ to server
- [Server] aggregates and sends $\{(\mathbf{x}_s, y_s)\}_{s \in \mathcal{S}_t}$ back to all clients, where $\mathcal{S}_t = \cup_{i \in [N]} \mathcal{S}_{t,i}$
- [Clients $\forall i$] compute and send $\{\mathbf{Z}_{\mathcal{N}_t(i), \mathcal{S}_t}^T, \mathbf{Z}_{\mathcal{N}_t(i), \mathcal{S}_t}, \mathbf{Z}_{\mathcal{N}_t(i), \mathcal{S}_t}^T \mathbf{y}_{\mathcal{N}_t(i)}\}$ to server
- [Server] aggregates $\sum_{i=1}^N \mathbf{Z}_{\mathcal{N}_t(i), \mathcal{S}_t}^T, \sum_{i=1}^N \mathbf{Z}_{\mathcal{N}_t(i), \mathcal{S}_t}, \sum_{i=1}^N \mathbf{Z}_{\mathcal{N}_t(i), \mathcal{S}_t}^T \mathbf{y}_{\mathcal{N}_t(i)}$ and sends it back
- [Clients $\forall i$] updates $\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}_t}^T, \mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}_t}$ and $\mathbf{Z}_{\mathcal{D}_t(i), \mathcal{S}_t}^T \mathbf{y}_{\mathcal{D}_t(i)}$; sets $\mathcal{D}_t(i) = \cup_{i=1}^N \mathcal{N}_t(i) = [t]$ and $t_{\text{last}} = t$

Theoretical Results

To attain near-optimal regret $R_T = O(\sqrt{NT}(\|\theta_*\| \sqrt{\gamma_{NT}} + \gamma_{NT}))$, our proposed solution requires $C_T = O(\gamma_{NT}^3 N^2)$ communication, where

- γ_{NT} is the maximum information gain, $\gamma_{NT} = d \log NT$ for linear kernel, $\gamma_{NT} = \log^{d+1} NT$ for Gaussian kernel
- under linear setting, it matches C_T of dedicated distributed linear bandit algorithms [Li and Wang, AISTATS' 22, He et al., NeurIPS' 22] up to $O(\log^2 NT)$

Experiment Results



References

- Yuanhao Wang, Jiachen Hu, Xiaoyu Chen, and Liwei Wang. (2019). "Distributed bandit learning: Near-optimal regret with efficient communication." In: International Conference on Learning Representations.
- Chuanhao Li and Hongning Wang (2022). "Asynchronous upper confidence bound algorithms for federated linear bandits." In: International Conference on Artificial Intelligence and Statistics.
- Jiafan He, Tianhao Wang, Yifei Min, Quanquan Gu (2022). "A Simple and Provably Efficient Algorithm for Asynchronous Federated Contextual Linear Bandits". In: Advances in Neural Information Processing Systems 36