(Private) Kernelized Bandits with **Distributed Biased Feedback** Fengjiao Li¹ Xingyu Zhou² Bo Ji¹

¹Virginia Tech ²Wayne State University

















Bayesian Optimization (BO)

• *f* lies in Reproducing Kernel Hilbert Space (**RKHS**)

- RKHS approximates **any** continuous function
- Include linear function as a special case
- Bounded RKHS norm, i.e., **smoothness**















Intuition&Fact Box

- Achieve $\tilde{O}(\gamma_T \sqrt{T})$ regret bound
- γ_T kernel-dependent, intrinsic dimensionality
 - Linear kernel: $O(d \log T)$, recover linear bandits
 - Gaussian kernel (SE): $\log^{d+1}(T)$
- Time complexity $O(T^3)$
 - Each GP update $O(t^2)$ using rank-1 update





What if... f is an expectation





15

Examples Policy making, cellular config and more...



Policy maximizes expected user satisfaction



How to apply BO (kernel bandits) in this case?







Base station config optimizes expected performance



Option I: Aggregate All Impractical...



How to apply BO (kernel bandits) in this case? $f(x) = \mathbb{E}_{u}F(x, u)$ Base station config optimizes expected performance X_t $f_1(x_t) + \eta$ 13 3 5 5 Large amount of comm.



Option II: Sample One GP-UCB/TS with a larger noise



Policy maximizes expected user satisfaction





- This is similar to SGD for Stochastic Optimization
 - One-pass over the dataset to have unbiased sar







Base station config optimizes expected performance

		Ì
nple		

Option II: Sample One GP-UCB/TS with larger noise





Three Limitations

• Poor scalability due to $O(T^3)$ time complexity



Intuition&Fact Box

• This is similar to SGD for Stoc. Optimization



Contribution



Policy maximizes expected user satisfaction



- Propose a distributed **phase-elimination-type** BO algorithm with **user sampling** and **batching**
- Involve with only $O(T^{\alpha})$ unique users, $\alpha \in (0,1]$ —user sampling parameter 2.
- 3.
- 4.











Policy maximizes expected user satisfaction

Phase-elimination

- Doubling phase length
- Eliminate poor actions/arms at end of each phase
- Can be easily extended to infinite domain setting







Policy maximizes expected user satisfaction



User Sampling

- Each phase samples a new set of users
- Those users are fixed during the phase
- Each phase l, $|U_l| = 2^{\alpha l}$, $\alpha \in (0,1]$





Policy maximizes expected user satisfaction



Action Selection

• Rarely switching via batching





Policy maximizes expected user satisfaction



Action Selection

- Rarely switching via batching
- Maximum variance reduction





Policy maximizes expected user satisfaction



Action Selection

- Rarely switching via batching
- Max-variance reduction principle







Theoretical Results

Performance Guarantees Non-private version

Theorem 1

There exists proper parameter choices of our algorithm such that it enjoys

- **1. Regret** $\tilde{O}(T^{1-\alpha/2} + \sqrt{\gamma_T T})$ - thanks to phase-elimination
- 2. Total number of users $O(T^{\alpha})$ - thanks to user reuse
- 3. Computation complexity $O(\gamma_T T^{\alpha})$ thanks to batching and merging
- 4. Communication cost $O(\gamma_T T^{\alpha})$

- thanks to batching and merging

- 1. Most related work is our prior work on linear bandit [LZJ'22] our results include it as a special case
- 2. Experimental design for action selection in kernel bandits [CJK21] our selection is much simpler without additional estimator
- 3. Zero-order non-convex stochastic opt.[BG18] our metric is total regret rather than convergence rate

Intuition&Fact Box

- GP-UCB achieves $\tilde{O}(\gamma_T \sqrt{T})$ regret bound
- γ_T kernel-dependent, intrinsic dimensionality
 - Linear kernel: $O(d \log T)$, recover linear bandits
 - Gaussian kernel (SE): $\log^{d+1}(T)$
- Time complexity $O(T^3)$ for GP-UCB

Compare with prior work

Privacy Protection

Differential Privacy

Differential Privacy 101

Definition. If for any two neighboring datasets D and D', and any outcome E $\mathbb{P}(M(D) \in E) \leq e^{e}\mathbb{P}(M(D') \in E) + \delta$ Then, M satisfies (e, δ) -DP - DP means that outputs are "close" in

probability on two neighboring datasets

Key components:

- 1. What are the neighboring datasets?— the identity for protection
- 2. What are the outputs?
 - the view of adversary

Key properties:

- 1. Composition, privacy loss adds up
- 2. Post-processing, immune to further processing if data is not touched

Differential Privacy Shuffle model

(î î)

Neighboring Dataset & Output

- Two user sequences differing in one
- All **inputs** of the central agent

100

Differential Privacy 301

- Shuffle DP leverages additional randomness in the shuffler to amplify privacy
- Same untrusted agent as local model
- Huge improvement of utility

Policy maximizes expected user satisfaction

3. Keep our privacy noise in control

Thank you!