

Follow-ups Also Matter: Improving Contextual Bandits via Post-serving Contexts

Chaoqi Wang¹, Ziyu Ye¹, Zhe Feng²,
Ashwinkumar Badanidiyuru³, Haifeng Xu¹

The University of Chicago¹
Google Research²
Google³

NeurIPS 2023

Background

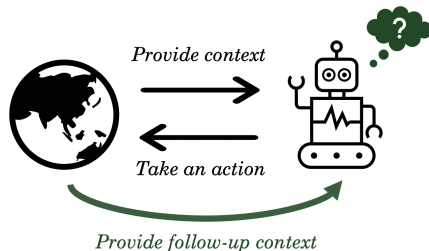


Figure 1: Illustration of learning with post-serving contexts.

- **Motivation:** Post-serving contexts are prevalent in recommender systems.
- **Challenges:** Classical bandit algorithms often fall short in such scenarios.
- **Research question:** How to effectively utilize **post-serving information** in **linear contextual bandits**?

Problem Setup and Notations

- ▶ **Problem Setup:** Each time $t = 1, 2, \dots, T$:
 - ▶ The learner observes the context \mathbf{x}_t .
 - ▶ The learner selects an arm $a_t \in [K]$.
 - ▶ The learner observes the reward r_{t,a_t} .
 - ▶ **The learner observes the post-serving context \mathbf{z}_t .**
- ▶ **Notations:**
 - ▶ Actions space: $\mathcal{A} = [K]$.
 - ▶ Pre-serving context: $\mathbf{x} \in \mathbb{R}^{d_x}$; post-serving context: $\mathbf{z} \in \mathbb{R}^{d_z}$.
 - ▶ $\mathbf{z} = \phi^*(\mathbf{x}_t) + \epsilon_t$, and $\phi^*(\mathbf{x}) = \mathbb{E}[\mathbf{z} \mid \mathbf{x}]$
 - ▶ Reward function:
 - ▶ $r_a(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \boldsymbol{\theta}_a^* + \mathbf{z}^\top \boldsymbol{\beta}_a^* + \eta$, where η is R_η -sub-Gaussian.
 - ▶ Matrix representation:
 - ▶ $\mathbf{X}_t = \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^\top + \lambda \mathbf{I}$ and $\mathbf{Z}_t = \sum_{s=1}^t \mathbf{z}_s \mathbf{z}_s^\top + \lambda \mathbf{I}$.
 - ▶ Norm restrictions:
 - ▶ $\forall a \in \mathcal{A}, \|\boldsymbol{\theta}_a^*\|_2 \leq 1, \|\boldsymbol{\beta}_a^*\|_2 \leq 1; \|\mathbf{x}\|_2 \leq L_x, \|\mathbf{z}\|_2 \leq L_z$.

Our Contributions

- ▶ **New framework:**
 - ▶ Proposed a novel family of contextual bandits with post-serving contexts.
- ▶ **Enhanced lemma:**
 - ▶ Introduced the Generalized Elliptical Potential Lemma (EPL).
- ▶ **Algorithm and theory:**
 - ▶ Designed poLinUCB with a regret bound of $\tilde{\mathcal{O}}(T^{1-\alpha} d_u^\alpha + d_u \sqrt{TK})$.
- ▶ **Empirical validation:**
 - ▶ Achieved SOTA performance on synthetic and real-world datasets.

Assumption: Generalized Learnability of $\phi^*(\cdot)$

Learnability Assumption

There exists an algorithm that, given t pairs of examples $\{(\mathbf{x}_s, \mathbf{z}_s)\}_{s=1}^t$ with arbitrarily chosen \mathbf{x}_s 's, outputs an estimated function of $\phi^* : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_z}$ such that for any $\mathbf{x} \in \mathbb{R}^{d_x}$, the following holds with probability at least $1 - \delta$,

$$e_t^\delta := \left\| \hat{\phi}_t(\mathbf{x}) - \phi^*(\mathbf{x}) \right\|_2 \leq C_0 \cdot \left(\|\mathbf{x}\|_{\mathbf{X}_t^{-1}}^2 \right)^\alpha \cdot \log(t/\delta),$$

where $\alpha \in (0, 1/2]$ and C_0 is some universal constant.

- ▶ The larger the value of α , the faster the learning rate for $\phi^*(\cdot)$.
- ▶ The regret of our algorithm is tied to $O(T^{1-\alpha})$.
- ▶ For linear functions, $\alpha = 1/2$ is optimal, yielding a regret of $O(\sqrt{T})$.

Why Natural Attempts May be Inadequate?

- ▶ Similar to [Wang et al., 2016]¹, a natural idea is to fit $\hat{\phi}(\cdot)$, and obtain the parameter estimate by solving:

$$\ell_t(\boldsymbol{\theta}_a, \boldsymbol{\beta}_a) = \sum_{s \in [t]: a_s = a} \left(r_{s,a} - \mathbf{x}_t^\top \boldsymbol{\theta}_a - \hat{\phi}_s(\mathbf{x}_s)^\top \boldsymbol{\beta}_a \right)^2 + \lambda (\|\boldsymbol{\theta}_a\|_2^2 + \|\boldsymbol{\beta}_a\|_2^2).$$

- ▶ The regret can be $\tilde{O}(T^{3/4})$ when initialized away from the global optimum.
- ▶ We propose to get the parameter estimate by solving:

$$\ell_t(\boldsymbol{\theta}_a, \boldsymbol{\beta}_a) = \sum_{s \in [t]: a_s = a} \left(r_{s,a} - \mathbf{x}_s^\top \boldsymbol{\theta}_a - \mathbf{z}_s^\top \boldsymbol{\beta}_a \right)^2 + \lambda (\|\boldsymbol{\theta}_a\|_2^2 + \|\boldsymbol{\beta}_a\|_2^2).$$

- ▶ This requires modification over the original Elliptical Potential Lemma (EPL) to accommodate noise in contexts during learning.

¹Huazheng Wang, Qingyun Wu, and Hongning Wang. “Learning Hidden Features for Contextual Bandits”. In: *CIKM*. 2016, pp. 1633–1642.

The Proposed Lemma: Generalized EPL

Generalized Elliptical Potential Lemma²

Suppose (1) $\mathbf{X}_0 \in \mathbb{R}^{d \times d}$ is any positive definite matrix; (2) $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^d$ and $\max_t \|\mathbf{x}_t\| \leq L_x$; (3) $\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_T \in \mathbb{R}^d$ are independent bounded zero-mean noises satisfying $\max_t \|\boldsymbol{\epsilon}_t\| \leq L_\epsilon$ and $\mathbb{E}[\boldsymbol{\epsilon}_t \boldsymbol{\epsilon}_t^\top] \succcurlyeq \sigma_\epsilon^2 \mathbf{I}$; and (4) $\widetilde{\mathbf{X}}_t$ is defined as:

$$\widetilde{\mathbf{X}}_t = \mathbf{X}_0 + \sum_{s=1}^t (\mathbf{x}_s + \boldsymbol{\epsilon}_s)(\mathbf{x}_s + \boldsymbol{\epsilon}_s)^\top \in \mathbb{R}^{d \times d}.$$

For any $p \in [0, 1]$, the following inequality holds with probability at least $1 - \delta$,

$$\sum_{t=1}^T \left(1 \wedge \|\mathbf{x}_t\|_{\widetilde{\mathbf{X}}_{t-1}^{-1}}^2 \right)^p \leq 2^p T^{1-p} \log^p \left(\frac{\det \mathbf{X}_T}{\det \mathbf{X}_0} \right) + \frac{8L_\epsilon^2(L_\epsilon + L_x)^2}{\sigma_\epsilon^4} \log \left(\frac{32dL_\epsilon^2(L_\epsilon + L_x)^2}{\delta\sigma_\epsilon^4} \right).$$

²The original EPL corresponds to the specific case of $p = 1$.

The Proposed Algorithm: poLinUCB

Algorithm 1 poLinUCB (Linear UCB with post-serving contexts)

- 1: **for** $t = 0, 1, \dots, T$ **do**
- 2: Receive the pre-serving context \mathbf{x}_t .
- 3: Compute the optimistic parameters by maximizing the UCB objective:

$$\left(a_t, \tilde{\phi}_t(\mathbf{x}_t), \tilde{\mathbf{w}}_t\right) = \arg \max_{(a, \phi, \mathbf{w}_a) \in [K] \times \mathcal{C}_{t-1}(\hat{\phi}_{t-1}, \mathbf{x}_t) \times \mathcal{C}_{t-1}(\hat{\mathbf{w}}_{t-1}, a)} \left[\begin{array}{c} \mathbf{x}_t \\ \phi(\mathbf{x}_t) \end{array} \right]^\top \mathbf{w}_a.$$

- 4: Play the arm a_t and receive the post-serving context \mathbf{z}_t and the reward r_{t,a_t} .
- 5: Compute $\hat{\mathbf{w}}_{t,a}$ for each $a \in \mathcal{A}$ using:

$$\ell_t(\boldsymbol{\theta}_a, \beta_a) = \sum_{s \in [t]: a_s = a} (r_{s,a} - \mathbf{x}_s^\top \boldsymbol{\theta}_a - \mathbf{z}_s^\top \beta_a)^2 + \lambda \left(\|\boldsymbol{\theta}_a\|_2^2 + \|\beta_a\|_2^2 \right).$$

- 6: Compute the estimated post-serving context generating function $\hat{\phi}_t(\cdot)$ using ERM.
 - 7: Update confidence sets $\mathcal{C}_t(\hat{\mathbf{w}}_{t,a})$ and $\mathcal{C}_t(\hat{\phi}_t, \mathbf{x}_t)$ for each a .
 - 8: **end for**
-

Regret Analysis

Settings	Ours
Ours	$\tilde{O}\left(T^{1-\alpha}d_u^\alpha + d_u\sqrt{TK}\right)$
Action-dependent contexts	$\tilde{O}\left(T^{1-\alpha}d_u^\alpha\sqrt{K} + d_u\sqrt{TK}\right)$
Same setting as in [Abbasi et al., 2011] ³	$\tilde{O}\left(T^{1-\alpha}d_u^\alpha + d_u\sqrt{T}\right)$

Table 1: Upper bound of regret of poLinUCB.

³Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. “Improved Algorithms for Linear Stochastic Bandits”. In: *Advances in neural information processing systems* 24 (2011).

Experimental Results: The Synthetic Dataset

- Our proposed poLinUCB consistently outperforms other strategies.
(Except for LinUCB (x and z) which equips with post-serving contexts in arm selection.)

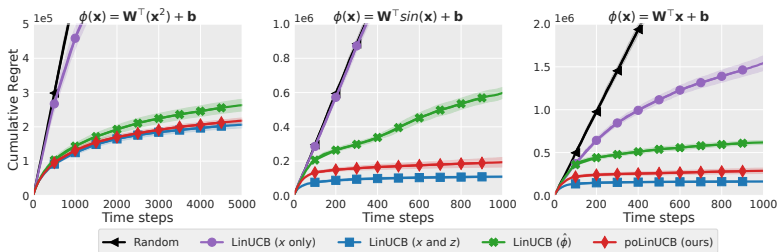


Figure 2: Algorithms' cumulative regrets in three synthetic environments. The shaded area denotes the standard error computed using 10 different random seeds.

Experimental Results: The MovieLens Dataset⁴

- Our proposed poLinUCB consistently outperforms other strategies.
(Except for LinUCB (x and z) which equips with post-serving contexts in arm selection.)

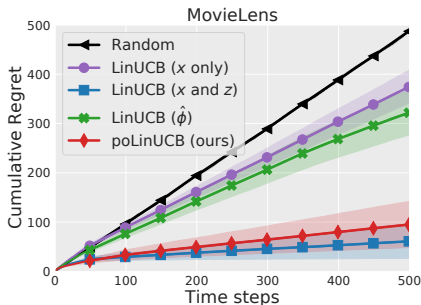


Figure 3: Algorithms' cumulative regrets in the MovieLens Dataset. The shaded area denotes the standard error computed using 10 different random seeds.

⁴F Maxwell Harper and Joseph A Konstan. "The MovieLens Datasets: History and Context". In: *Acm Transactions on Interactive Intelligent Systems* 5.4 (2015), pp. 1–19.

Thank you.

Please refer to our paper for more information:

<https://arxiv.org/abs/2309.13896>.