

Near-Optimal Representation Learning for Linear Bandits and Linear RL

Jiachen Hu ¹ Xiaoyu Chen ¹ Chi Jin ² Lihong Li ³ Liwei Wang ¹

¹Peking University ²Princeton University ³Amazon

June 20, 2021

Multi-Task Low Rank Linear Bandits

- The agent chooses arms for M linear bandits instance concurrently
 - ▶ At each step, he pulls M arms and then receives M rewards
 - ▶ Linear rewards $r_{t,i} = \langle a_{t,i}, \theta_i^* \rangle + \text{noise}$
- The **low rank** property: There exists $k \ll d$ such that

$$\text{rank} \left(\begin{bmatrix} \theta_1^* & \theta_2^* & \cdots & \theta_M^* \end{bmatrix} \right) = k$$

$$\theta_i^* = B^* w_i^*, \forall i \in [M]$$

- $B^* \in \mathbb{R}^{d \times k}$: the shared k -dim representation. $w_i^* \in \mathbb{R}^k$: the task-dependent parameter

Multi-Task Low Rank Linear Bandits

- The agent chooses arms for M linear bandits instance concurrently
 - ▶ At each step, he pulls M arms and then receives M rewards
 - ▶ Linear rewards $r_{t,i} = \langle a_{t,i}, \theta_i^* \rangle + \text{noise}$
- The **low rank** property: There exists $k \ll d$ such that

$$\text{rank} \left(\begin{bmatrix} \theta_1^* & \theta_2^* & \cdots & \theta_M^* \end{bmatrix} \right) = k$$

$$\theta_i^* = B^* w_i^*, \forall i \in [M]$$

- $B^* \in \mathbb{R}^{d \times k}$: the shared k -dim representation. $w_i^* \in \mathbb{R}^k$: the task-dependent parameter
- **We can estimate B^* using the data from all tasks!**

Algorithm

Multi-task linear bandits

For step $t = 1, 2, \dots, T$:

- 1 Estimate $\{\theta_i^*\}_{i=1}^M$ using **constrained** least-square linear regression

$$\{\hat{\theta}_{t,i}\}_{i=1}^M = \{\hat{B}_t \hat{w}_{t,i}\}_{i=1}^M \stackrel{\text{def}}{=} \arg \min_{\|Bw_i\|_2 \leq 1} \sum_{\tau=1}^{t-1} \sum_{i=1}^M (\langle a_{\tau,i}, Bw_i \rangle - r_{\tau,i})^2$$

- 2 Construct **tighter** confidence set $\mathcal{C}_t \subset (\mathbb{R}^d)^M$ centered at $\{\hat{\theta}_{t,i}\}_{i=1}^M$ such that $\{\theta_i^*\}_{i=1}^M \in \mathcal{C}_t$ with high probability
- 3 Choose actions optimistically and observe reward $\{r_{t,i}\}_{i=1}^M$

$$\{a_{t,i}\}_{i=1}^M \leftarrow \arg \max_{a_i \in \mathcal{A}_{t,i}} \max_{\{\theta_i\}_{i=1}^M \in \mathcal{C}_t} \sum_{i=1}^M \langle a_i, \theta_i \rangle$$

Results

- Independent single-task algorithm for each instance

$$\tilde{O}\left(Md\sqrt{T}\right)$$

- Multi-task algorithm

$$\tilde{O}\left(M\sqrt{dkT} + d\sqrt{MkT}\right)$$

- Known lower bound

$$\Omega\left(Mk\sqrt{T} + d\sqrt{MkT}\right)$$

Multi-task Low Rank Episodic RL

- The multi-task LSVI condition: We have access to a **joint linear function class \mathcal{Q}** with **low inherent Bellman error**¹
- Linear parameters in \mathcal{Q} share a k -dim subspace
 - ▶ k is much smaller than the feature dimension d

¹See Zanette et al. Learning Near Optimal Policies with Low Inherent Bellman Error

Multi-task Low Rank Episodic RL

- The multi-task LSVI condition: We have access to a **joint linear function class \mathcal{Q}** with **low inherent Bellman error**¹
- Linear parameters in \mathcal{Q} share a k -dim subspace
 - ▶ k is much smaller than the feature dimension d
- Independent single-task regret

$$\tilde{O}\left(HMd\sqrt{T} + HMT\sqrt{d\mathcal{I}}\right)$$

- **Multi-task regret**

$$\tilde{O}\left(HM\sqrt{dkT} + Hd\sqrt{kMT} + HMT\sqrt{d\mathcal{I}}\right)$$

- **Lower bound**

$$\Omega\left(Mk\sqrt{HT} + d\sqrt{HkMT} + HMT\sqrt{d\mathcal{I}}\right)$$

¹See Zanette et al. Learning Near Optimal Policies with Low Inherent Bellman Error