

# MULTI-TASK LINEAR BANDITS

MARTA SOARE, OUAIS ALSHARIF, ALESSANDRO LAZARIC, JOELLE PINEAU

## MOTIVATION

- **Recommendation System:** users with similar features have similar preferences over different items.
- **Personalized Healthcare:** patients with similar symptoms and medical history have the similar reactions to treatments.
- **Games:** an agent might reduce the exploratory steps needed to discover an environment, by using the knowledge acquired on previous similar environments.

## PROBLEM SETTING

**Sequential multi-task linear bandit:** The learner faces a sequence of (unknown) linear bandit tasks  $(\theta_1, \theta_2, \dots, \theta_j, \dots)$ .

**Linear bandit task:**

- Set of arms  $\mathcal{X} \subseteq \mathbb{R}^d$ ,  $|\mathcal{X}| = K$  and  $\|x\|_2 \leq L$ ,  $\forall x \in \mathcal{X}$ .
- Reward model for task  $j$ :

$$r_j(x) = x^\top \theta_j + \eta$$

where  $\theta_j \in \mathbb{R}^d$  is **unknown** and  $\eta$  is an R sub-Gaussian noise.

- Cumulative regret wrt the optimal arm  $x_j^* = \arg \max_{x \in \mathcal{X}} x^\top \theta_j$

$$R_{n_j} = \sum_{s=1}^{n_j} (x_j^* - x_s)^\top \theta_j$$

**Task Similarity:** there exists  $\varepsilon > 0$  such that for any pair of tasks  $(\theta, \theta')$  we have  $\|\theta - \theta'\|_2 \leq \varepsilon$ .

**Better performance on a particular task can be achieved by leveraging information from different but similar tasks.**

## MAIN RESULT

**Theorem 1** Let  $\tilde{\theta}_{m+1,t}^\lambda$  be the multi-task regularized least-squares estimate. Then, for any  $\delta \geq 0$ , for any  $t \geq 1$ , with probability greater than  $1 - \delta$  it holds that:

$$\|x^\top (\tilde{\theta}_{m+1,t}^\lambda - \theta_{m+1})\| \leq \|x\| (\tilde{A}_{m+1,t}^\lambda)^{-1} \left( R \sqrt{d \log \left( \frac{\det(\tilde{A}_{m+1,t}^\lambda)^{1/2} \det(\lambda I)^{-1/2}}{\delta} \right)} + \lambda^{1/2} S \right) + x^\top \varepsilon = \tilde{B}_{m+1,t}(x).$$

## ESTIMATION ERROR

The estimation error of the MT least-squares estimate is upper-bounded by

$$\|x^\top (\tilde{\theta}_{m+1,t}^\lambda - \tilde{\theta}_{m+1,t}^*)\| \leq \|x\| (\tilde{A}_{m+1,t}^\lambda)^{-1} \left( R \sqrt{2 \log \left( \frac{\det(\tilde{A}_{m+1,t}^\lambda)^{1/2} \det(\lambda I)^{-1/2}}{\delta} \right)} + \lambda^{1/2} S \right)$$

## MULTI-TASK BIAS

Under the **task similarity assumption** the approximation error of the multi-task estimates is

$$\|\tilde{\theta}_{m+1,t}^* - \theta_{m+1}\| \leq \varepsilon, \text{ w.p. } 1 - \delta.$$

## TOOLS

**Single-task ordinary least-squares (OLS) estimate:** for any task  $j$ , the estimate of  $\theta_j$  after  $n$  rewards is

$$A_{j,n} = \sum_{t=1}^n x_t x_t^\top, \quad b_{j,n} = \sum_{t=1}^n x_t r_t, \quad \hat{\theta}_{j,n} = A_{j,n}^{-1} b_{j,n}$$

**Single-task regularized least-squares estimate:** for any task  $j$ , the estimate of  $\theta_j$  after  $n_j$  rewards is

$$A_{j,n}^\lambda = \left( \sum_{t=1}^n x_t x_t^\top + \lambda I \right), \quad b_{j,n} = \sum_{t=1}^n x_t r_t, \quad \hat{\theta}_{j,n}^\lambda = (A_{j,n}^\lambda)^{-1} b_{j,n}$$

**Single-task prediction error** Thm.2 in [1] with probability  $1 - \delta$

$$\|x^\top \theta_j - x^\top \hat{\theta}_{j,n}^\lambda\| \leq \|x\| (A_{j,n}^\lambda)^{-1} \left( R \sqrt{d \log \left( \frac{1 + nL^2/\lambda}{\delta} \right)} + \lambda^{1/2} \|\theta_j\| \right) = B_{j,n}(x).$$

## MULTI-TASK ESTIMATES

**Multi-task estimates:** use all the past samples to construct an estimate for the current task

$$\tilde{A}_{m+1,t} = \sum_{j=1}^m A_j + A_{m+1,t}; \quad \tilde{b}_{m+1,t} = \sum_{j=1}^m b_j + b_{m+1,t}$$

$$\tilde{\theta}_{m+1,t} = \tilde{A}_{m+1,t}^{-1} \tilde{b}_{m+1,t}$$

**Average target task:**

$$\mathbb{E}[\tilde{\theta}_{m+1,t}] = \tilde{\theta}_{m+1,t}^*$$

**Multi-task Regularized estimates:**

$$\tilde{A}_{m+1,t}^\lambda = \tilde{A}_{m+1,t} + \lambda I, \quad \tilde{\theta}_{m+1,t}^\lambda = (\tilde{A}_{m+1,t}^\lambda)^{-1} \cdot \tilde{b}_{m+1,t}$$

**Use at the same time single-task and multi-task estimates to construct upper confidence bounds  $B(x), \tilde{B}(x)$ .**

## MULTITASK-LINUCB

**Input:** budgets  $\{n_j\}_j$ , arms  $\mathcal{X} \subset \mathbb{R}^d$ , regularizer  $\lambda$

$j = 1$   
 $A = \lambda I_d, \tilde{b} = b = 0_d, \tilde{A}_j = \lambda I_d, \hat{\theta}_j = A^{-1} b$

**for**  $t = 1, \dots, n_j$  **do**

Choose:  $x_t = \arg \max_{x \in \mathcal{X}} (x_t^\top \hat{\theta}_j + B_{j,t}(x))$

Observe reward:  $r_t = x_t^\top \theta_j + \eta_t$

Update  $A, b$  and the estimate  $\hat{\theta}_j = A^{-1} b$

**end for**

**for**  $j = 2, \dots, m + 1$  **do**

$\tilde{A}_j = \tilde{A}_j + A - \lambda I_d, \tilde{b} = \tilde{b} + b, \tilde{\theta}_j = \tilde{A}_j^{-1} \tilde{b}$

$A = \lambda I_d, b = 0_d, \hat{\theta}_j = A^{-1} b$

**for**  $t = 1, \dots, n_j$  **do**

$x_t = \arg \max_{x \in \mathcal{X}} \min [x^\top \hat{\theta}_j + B_{j,t}(x); x^\top \tilde{\theta}_j + \tilde{B}_{j,t}(x)]$

Observe reward:  $r_t = x_t^\top \theta_j + \eta_t$

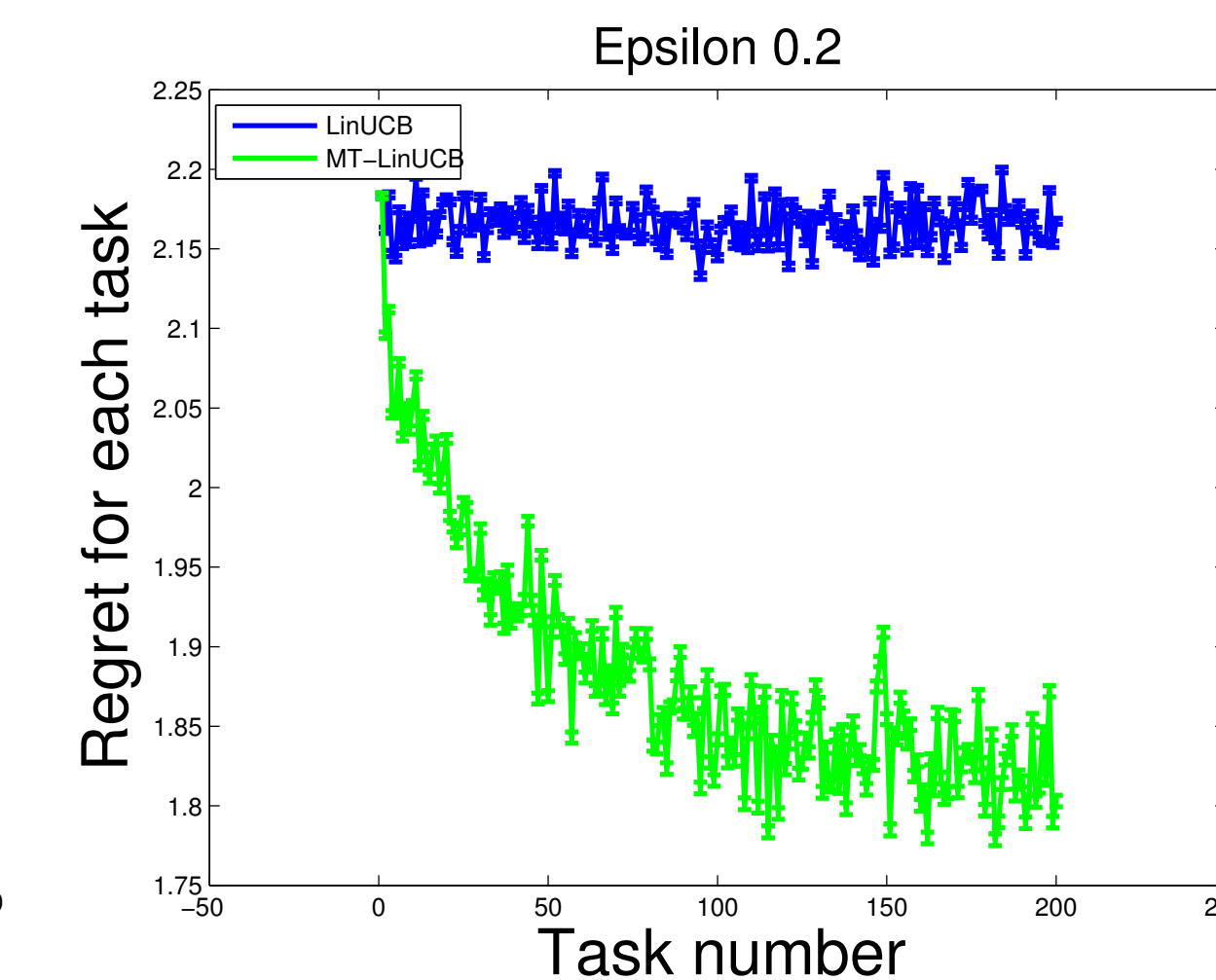
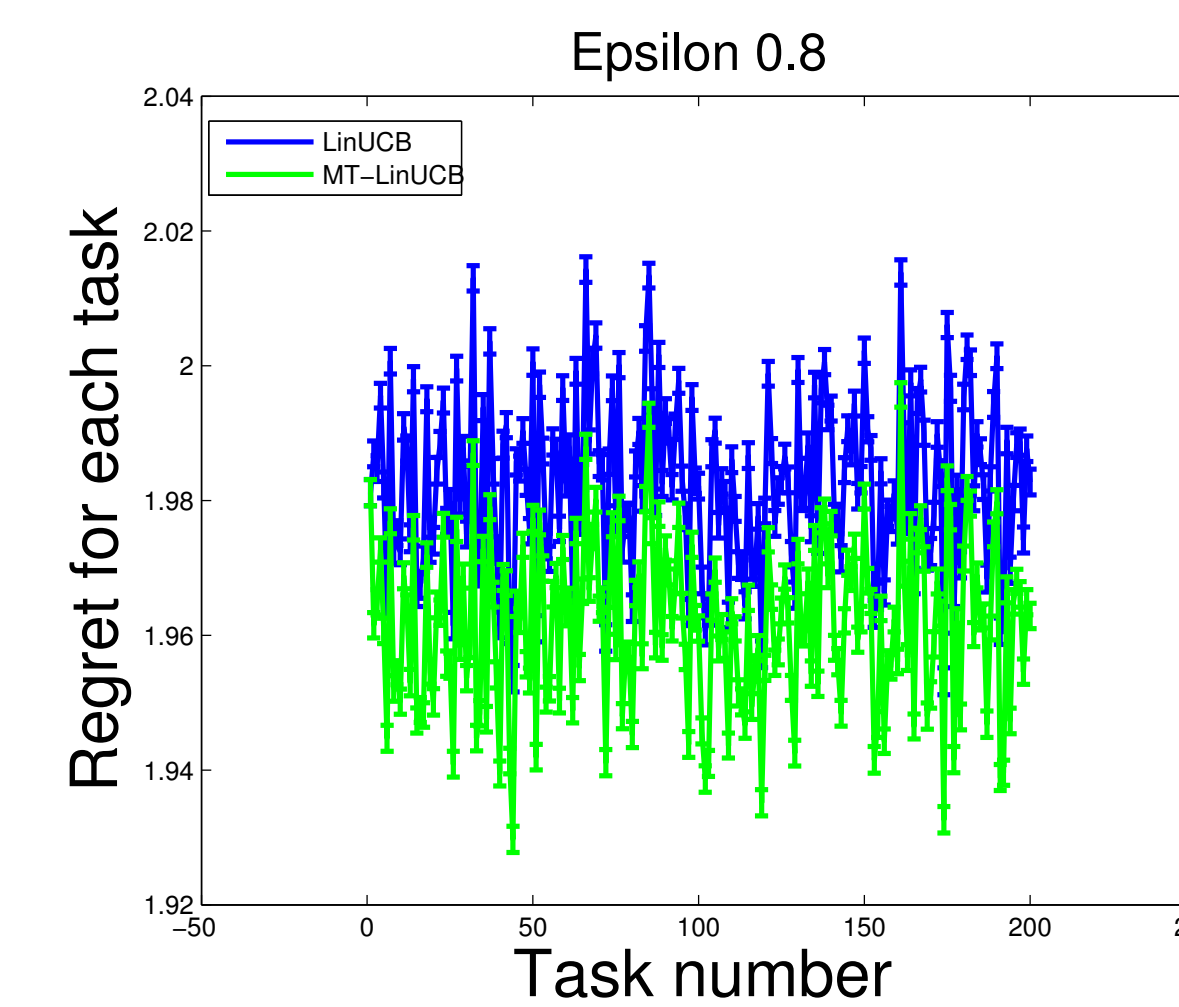
Update:  $A, b, \hat{\theta}_j, B_{j,t}, \tilde{A}_j, \tilde{b}, \tilde{\theta}_j, \tilde{B}_{j,t}$

**end for**

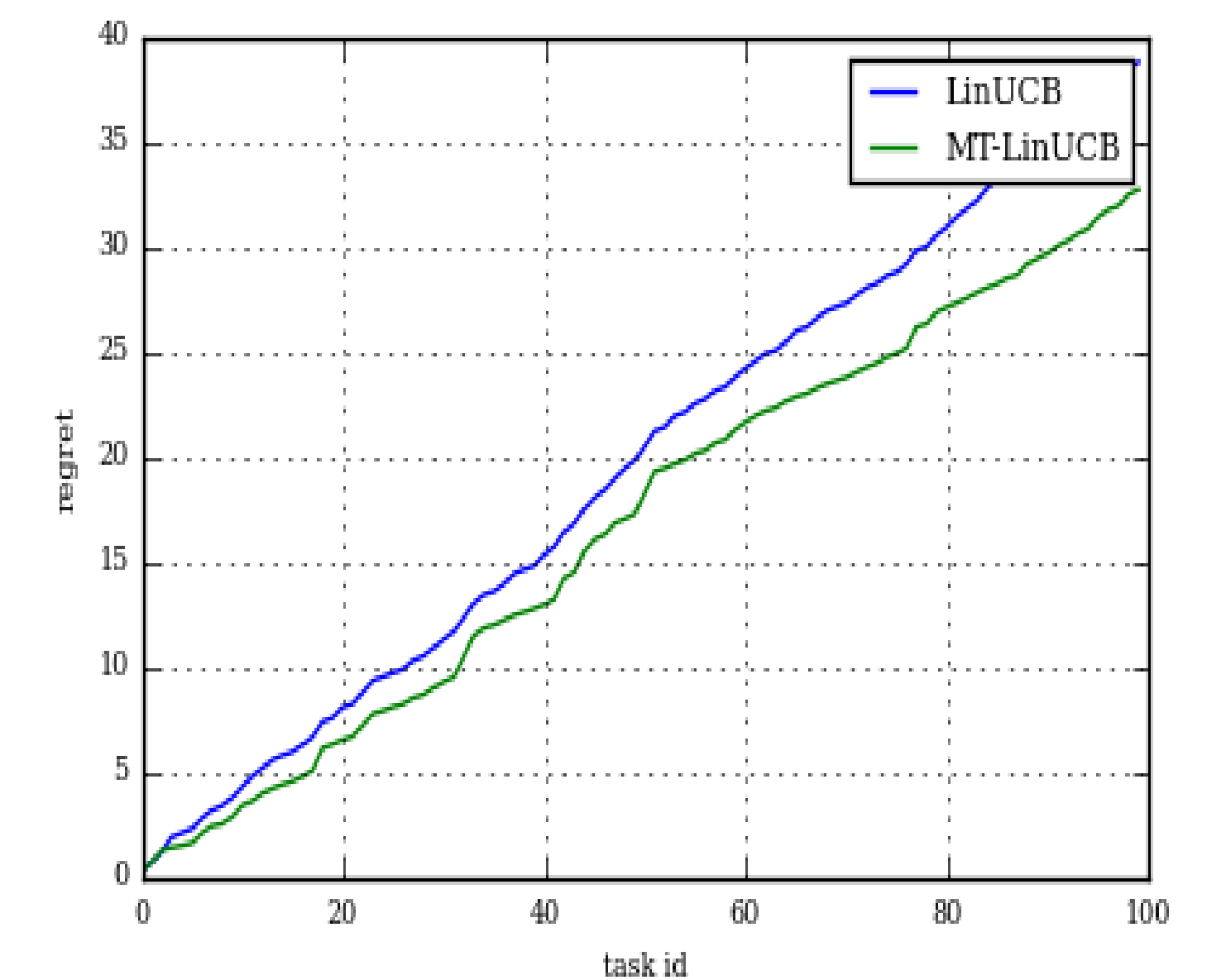
**end for**

## EXPERIMENTS

### Synthetic experiments



### Boston Housing dataset



### Setting

- 200 tasks with  $\theta_1, \dots, \theta_{200} \in \mathbb{R}^2$ , randomly generated
- $\max_{\theta} \|\theta\|_2 = 1.1 \cdot \sqrt{2}$ ,  $\|\varepsilon\|_2$  ranges in  $\{0.8, 0.6, 0.4, 0.2\} \cdot \sqrt{2}$
- 100 samples for each task

### Results

- MT-LINUCB is never worse than LINUCB.
- As the difference between tasks reduces, the advantage of MT-LINUCB becomes more evident.
- For  $\varepsilon = 0.2\sqrt{2}$ , the regret of MT-LINUCB decreases with every additional task, while the regret for LINUCB remains constant over time.

- Data points are scaled to have a norm of 1.
- Twenty different clusters, each cluster center is  $\theta_m^*$  for task  $m$ .
- For each task, we drew 10 sets of 3 arms each.

## REFERENCES

[1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances Neural Information Processing Systems* 24, 2011.

## ACKNOWLEDGEMENTS

